

\*\* Result [Patent] \*\* Format(P801) 01. Sep. 2003 1/ 1

Application no/date: 2000-526867[1998/12/21]  
Date of request for examination: [ ]  
Public disclosure no/date: 2002-500393[2002/01/08] **Translate**  
Examined publication no/date (old law): [ ]  
Registration no/date: [ ]  
Examined publication date (present law): [ ]  
PCT application no PCT/US98/27199  
PCT publication no/date WO/99/034291[1999/07/08]  
Applicant: AVID TECHNOL INC  
Inventor: PIITAAZU ERITSUKU SHII, RABINOITSUTSU SUTANREE, JIEIKOBUSU HAABAAT  
O AARU, JIRETSUTO RICHIIYAADO BEIKAA JIYUNIA, FUATSUSHIANO PIITAA JIEI  
IPC: G06F 13/10 ,340 G06F 3/06 ,305 G06F 3/06 ,540  
G06F 11/08 ,310 G06F 12/00 ,545 G06F 12/16 ,320  
H04N 7/173 ,610  
FI: G06F 11/08 ,310B G06F 13/10 ,340A G06F 12/16 ,320L  
G06F 3/06 ,305C G06F 3/06 ,540 G06F 12/00 ,545B H04N 7/173 ,610A  
F-term: 5B001AA00, AB01, AC01, AD03, 5B014EA04, EB04, FB03, FB04, GC07, GD23, GD32,  
GD33, HA09, 5B018GA01, HA11, MA11, MA15, 5B065BA01, EA03, EA12, EA19, EA24, EA31, EA35,  
ZA08, ZA15, 5B082DE05, HA01, 5C064BA07, BB06, BC18, BD02, BD13  
Expanded classification: 452, 446, 451, 453  
Fixed keyword: R011, R012, R131  
Citation:  
Title of invention: The process which in scalable and dependability are high, and transfer hi  
Abstract:

PURPOSE: Invocation, storage memory storing segment every each segment  
of chosen data are chosen, reading, series become in each segment  
of the data which called, transfer of independent high bandwidth data  
stream of high mass of dependability is enabled by supplying data  
in an application when data is received from the storage memory which  
syoteishi.

CONSTITUTION: Capture system makes segment table 90A. The pictorial image  
index which a map makes each pictorial image in off set to a data  
stream to take in is made. For example, the data which became index  
can support field or a frame. In addition, Capture system acquires list  
of available storage memory. Furthermore, when an application program  
to carry out on a client called for data, when each data segment is  
stored on at least two storage memory, claim seems to be satisfied,  
and storage memory is chosen.

( Machine Translation )

---

(19) 日本国特許庁 (J P)

(12) 公表特許公報 (A)

(11) 特許出願公表番号  
特表2002-500393  
(P2002-500393A)

(43) 公表日 平成14年1月8日(2002.1.8)

(51) Int.Cl. <sup>7</sup>	識別記号	F I	テマコード <sup>*</sup> (参考)
G 0 6 F 13/10	3 4 0	G 0 6 F 13/10	3 4 0 A 5 B 0 0 1
3/06	3 0 5	3/06	3 0 5 C 5 B 0 1 4
	5 4 0		5 4 0 5 B 0 1 8
11/08	3 1 0	11/08	3 1 0 B 5 B 0 6 5
12/00	5 4 5	12/00	5 4 5 B 5 B 0 8 2
審査請求 未請求 予備審査請求 有 (全 124 頁) 最終頁に続く			

(21) 出願番号 特願2000-526867(P2000-526867)  
(86) (22) 出願日 平成10年12月21日(1998.12.21)  
(85) 翻訳文提出日 平成12年6月26日(2000.6.26)  
(86) 国際出願番号 PCT/US98/27199  
(87) 国際公開番号 WO99/34291  
(87) 国際公開日 平成11年7月8日(1999.7.8)  
(31) 優先権主張番号 08/997, 769  
(32) 優先日 平成9年12月24日(1997.12.24)  
(33) 優先権主張国 米国 (US)  
(31) 優先権主張番号 09/006, 070  
(32) 優先日 平成10年1月12日(1998.1.12)  
(33) 優先権主張国 米国 (US)

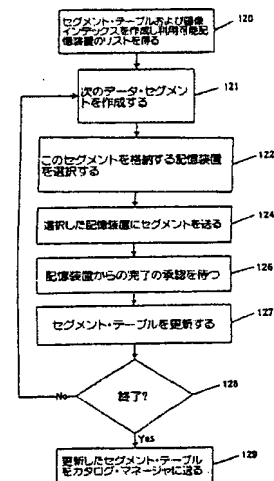
(71) 出願人 アヴィッド・テクノロジー・インコーポレ  
ーテッド  
AVID TECHNOLOGY, IN  
C.  
アメリカ合衆国マサチューセッツ州01876,  
テュークスバリー, ワン・パーク・ウエス  
ト, メトロポリタン・テクノロジー・パー  
ク  
(72) 発明者 ビーターズ, エリック・シー  
アメリカ合衆国マサチューセッツ州01741,  
カーライル, カールトン・ロード 80  
(74) 代理人 弁理士 社本 一夫 (外4名)

最終頁に続く

(54) 【発明の名称】 コンピュータ・システムおよび多数の記憶装置および多数のアプリケーション間でスケーラブル  
にかつ信頼性高く多数の高帯域データ・ストリームを転送するプロセス

#### (57) 【要約】

多数のアプリケーションは、コンピュータ・ネットワークを通じて多数の記憶装置からデータを要求する。データをセグメントに分割し、メディア・データの他のセグメントを格納した記憶装置には独立して、各セグメントを数個の記憶装置の1つにランダムに分散する。各セグメントに対応する冗長情報も、記憶装置間でランダムに分散する。セグメントに対する冗長情報は、セグメントのコピーとすることができ、各セグメントを少なくとも2つの記憶装置上に格納する。また、冗長情報は、2つ以上のセグメントに基づくことも可能である。このデータ・セグメントおよび対応する冗長情報のランダム分散により、スケーラビリティおよび信頼性双方の向上を図る。記憶装置が故障した場合、その負荷を均等に残りの記憶装置に分散し、その失われたデータを、冗長情報によって復元することができる。アプリケーションが選択したデータ・セグメントを要求した場合、最も短い要求キューを有する記憶装置がこの要求を処理することができる。多数のアプリケーションによって多数の記憶装置上に与えられる負荷のランダムな変動は、ほぼ等しく全



## 【特許請求の範囲】

【請求項1】 データを格納するための複数の記憶装置からなり、前記記憶装置上に格納したデータ・セグメントおよび対応する冗長情報を前記複数の記憶装置間でランダムに分散する、分散データ記憶システム。

【請求項2】 セグメントに対応する前記情報性情報は、当該セグメントのコピーである、請求項1記載の分散データ記憶システム。

【請求項3】 各セグメントの各コピーを前記記憶装置の異なるものに格納する請求項2記載の分散記憶システム。

【請求項4】 各セグメントの各コピーは、前記記憶装置の相対的仕様の関数として定義した確率分布にしたがって、前記複数の記憶装置の1つに割り当てられる請求項3記載の分散データ記憶システム。

【請求項5】 更に、コンピュータ読み取り可能ロジックを格納し、データ・セグメントの指示を用いてコンピュータによるアクセスが可能なセグメント・テーブルを定義するコンピュータ読み取り可能媒体を備え、前記セグメントおよび対応する冗長情報を格納した前記複数の記憶装置から、該記憶装置の指示を検索する請求項1記載の分散データ記憶システム。

【請求項6】 前記複数の記憶装置が、  
コンピュータ・ネットワークに接続した第1記憶装置と、  
前記コンピュータ・ネットワークに接続した第2記憶装置と、  
前記コンピュータ・ネットワークに接続した第3記憶装置と、  
を備える請求項1記載の分散データ記憶システム。

【請求項7】 セグメントに対応する前記冗長情報は、2つ以上のセグメントに基づく請求項1記載の分散データ記憶システム。

【請求項8】 前記2つ以上のセグメントおよび前記冗長情報を前記記憶装置の異なるものに格納する請求項7記載の分散データ記憶システム。

【請求項9】 コンピュータ用ファイル・システムであって、前記コンピュータが、当該コンピュータ上で実行するアプリケーションからの要求に回答してコンピュータ・ネットワークを通じて独立したリモート記憶装置にアクセスし、前記記憶装置上に格納してあるデータを読み出すことを可能とし、前記データの

セグメントおよび対応する冗長情報を前記複数の記憶装置間でランダムに分散するファイル・システムにおいて、

前記要求に応答してデータを読み出し、前記選択したデータの各セグメント毎に、前記セグメントを格納する記憶装置を特定する手段と、

前記要求したデータの各セグメントを、当該セグメントに対して特定した記憶装置から、読み出す手段と、

前記特定した記憶装置からデータを受信したとき、前記アプリケーションに前記データを供給する手段と、

を備えるファイル・システム。

【請求項10】 セグメントに対応する前記冗長情報は前記セグメントのコピーであり、前記セグメントを記憶する前記記憶装置を特定する前記手段は、前記選択したデータの各セグメント毎に、当該セグメントを格納する記憶装置の1つを選択する手段を含む請求項9記載のファイル・システム。

【請求項11】 前記記憶装置の1つを選択する手段は、前記複数の記憶装置上の要求の負荷が実質的に均衡するように記憶装置を選択することを特徴とする請求項10記載のファイル・システム。

【請求項12】 前記選択手段は、前記セグメントのためのどの記憶装置が、前記要求に応じるための推定時間が最も短いかについての推定に基づいて、前記セグメントのための前記記憶装置を選択する請求項11記載のファイル・システム。

【請求項13】 前記選択手段が、  
前記ファイル・システムにおいて、

前記記憶装置の1つからデータを要求し、推定時間を示す手段と、

前記第1記憶装置が前記要求を拒絶した場合、前記記憶装置の他のものからデータを要求し、推定時間を示す手段と、

前記第2記憶装置が前記要求を拒絶した場合、前記第1記憶装置から前記データを要求する手段と、

を含み、

各記憶装置において、

前記要求が前記推定時間以内に前記記憶装置によって対応することができない場合、データの要求を拒絶する手段と、

前記要求が前記推定時間以内に前記記憶装置によって対応することができる場合、データの要求を受け入れる手段と、

を含む請求項12記載のファイル・システム。

【請求項14】 各セグメントを読み出す前記手段は、前記選択した記憶装置からの前記データの転送をスケジュールし、前記記憶装置が効率的にデータ転送を行なうようにする手段を備える請求項11記載のファイル・システム。

【請求項15】 転送をスケジュールする前記手段は、

前記ファイル・システムにおいて、

前記選択した記憶装置から前記データの転送を要求し、待ち時間を示す手段と

前記選択した記憶装置が前記データを転送する前記要求を拒絶した場合、別の記憶装置から前記データを要求する手段と、

を含み、

前記記憶装置において、

前記データが前記指示した待ち時間までに前記記憶装置から転送することができない場合、データ転送要求を拒絶する手段と、

前記選択した記憶装置が前記待ち時間以内に前記データを転送することができる場合、前記データを転送する手段と、

を含む請求項14記載のファイル・システム。

【請求項16】 各セグメントを読み出す前記手段は、前記選択した記憶装置からの前記データの転送をスケジュールし、前記記憶装置が効率的にデータ転送を行なうようにする手段を備える請求項9記載のファイル・システム。

【請求項17】 各セグメントを読み出す前記手段は、前記選択した記憶装置からの前記データの転送をスケジュールし、前記記憶装置が効率的にデータ転送を行なうようにする手段を備える請求項7記載のファイル・システム。

【請求項18】 転送をスケジュールする前記手段は、

前記ファイル・システムにおいて、

前記選択した記憶装置から前記データの転送を要求し、待ち時間を示す手段と

前記選択した記憶装置が前記データを転送する前記要求を拒絶した場合、別の記憶装置から前記データを要求する手段と、

を含み、

前記記憶装置において、

前記データが前記指示した待ち時間までに前記記憶装置から転送することができない場合、データ転送要求を拒絶する手段と、

前記選択した記憶装置が前記待ち時間以内に前記データを転送することができる場合、前記データを転送する手段と、

を含む請求項17記載のファイル・システム。

【請求項19】 前記データを複数のセグメントに分割し、各セグメントをコピーし、各セグメントの各コピーを前記記憶装置の異なるものに格納する請求項10記載のファイル・システム。

【請求項20】 各セグメントの各コピーを、前記記憶装置の相対的仕様の関数として定義した確率分布にしたがって、前記複数の記憶装置の1つに割り当てる請求項19記載のファイル・システム。

【請求項21】 更に、コンピュータ読み取り可能ロジックを格納し、データ・セグメントの指示を用いてコンピュータによるアクセスが可能なセグメント・テーブルを定義するコンピュータ読み取り可能媒体を備え、前記セグメントおよび対応する冗長情報を格納した前記複数の記憶装置から、該記憶装置の指示を検索する請求項9記載のファイル・システム。

【請求項22】 前記複数の記憶装置が、  
前記コンピュータ・ネットワークに接続した第1記憶装置と、  
前記コンピュータ・ネットワークに接続した第2記憶装置と、  
前記コンピュータ・ネットワークに接続した第3記憶装置と、  
を備える請求項9記載のファイル・システム。

【請求項23】 セグメントに対応する前記冗長情報は、2つ以上のセグメントに基づく請求項9記載の分散データ記憶システム。

【請求項24】 前記冗長情報および該冗長情報の基となる前記2つ以上のセグメントは、各々前記記憶装置の異なる1つに記憶されている請求項23記載のファイル・システム。

【請求項25】 コンピュータ用ファイル・システムであって、前記コンピュータが、当該コンピュータ上で実行するアプリケーションからの、要求に応答してコンピュータ・ネットワークを通じて独立したリモート記憶装置にアクセスし、前記記憶装置上にデータを格納することを可能とするファイル・システムにおいて、

前記データを格納する前記要求に応答して、前記データを複数のセグメントに分割する手段と、

各セグメントおよび該セグメントに対応する冗長情報を前記複数の記憶装置間でランダムに分散する手段と、

前記データを格納したか否かについて前記アプリケーションに確認する手段と、  
を備えるファイル・システム。

【請求項26】 セグメントに対応する前記冗長情報は当該セグメントのコピーである請求項25記載のファイル・システム。

【請求項27】 前記ランダムに分散する手段は、  
各セグメント毎に、ランダムにかつ他のセグメントのために選択した記憶装置とは独立して、少なくとも2つの前記記憶装置を選択する手段と、

前記選択した記憶手段に、各セグメントに対するデータを格納するように要求する手段と、  
を備える請求項26記載のファイル・システム。

【請求項28】 前記選択する手段は、前記記憶装置の部分集合を選択する手段と、前記選択した部分集合内にある前記記憶装置の中から少なくとも2つの前記記憶装置を選択する手段とを含む請求項27記載のファイル・システム。

【請求項29】 各セグメントの各コピーは、前記記憶装置の相対的仕様の関数として定義した確率分布にしたがって、前記複数の記憶装置の1つに割り当てられる請求項27記載のファイル・システム。

【請求項30】 更に、コンピュータ読み取り可能ロジックを格納し、データ・セグメントの指示を用いてコンピュータによるアクセスが可能なセグメント・テーブルを定義するコンピュータ読み取り可能媒体を備え、前記セグメントおよび対応する冗長情報を格納した前記複数の記憶装置から、該記憶装置の指示を検索する請求項25記載のファイル・システム。

【請求項31】 前記複数の記憶装置が、  
前記コンピュータ・ネットワークに接続した第1記憶装置と、  
前記コンピュータ・ネットワークに接続した第2記憶装置と、  
前記コンピュータ・ネットワークに接続した第3記憶装置と、  
を備える請求項25記載のファイル・システム。

【請求項32】 セグメントに対応する前記冗長情報は、2つ以上のセグメントに基づく請求項25記載の分散データ記憶システム。

【請求項33】 前記冗長情報および該冗長情報の基となる前記2つ以上のセグメントは、各々前記記憶装置の異なる1つに記憶されている請求項32記載のファイル・システム。

【請求項34】 データを格納するための複数の記憶装置を備えた分散データ記憶システムにおいてデータを復元するプロセスであって、前記記憶装置上に格納した前記データのセグメントおよび対応する冗長情報は、前記複数の記憶装置間にランダムに分散しており、前記プロセスは、前記記憶装置の1つの故障が検出されたときに実行され、

前記故障した記憶装置上に格納してあったコピーのセグメントを特定するステップと、

前記特定したセグメントに対応する前記冗長情報を格納した記憶装置を特定するステップと、

前記複数の記憶装置間において、前記特定した冗長情報に応じて前記特定したセグメントをランダムに分散するステップと、  
からなるプロセス。

【請求項35】 セグメントに対応する前記冗長情報は当該セグメントのコピーであり、



記憶装置を特定する前記ステップは、前記特定したセグメントの別のコピーがストレージであった記憶装置を特定するステップを含み、

ランダムに分散する前記ステップは、前記特定した記憶装置からの前記特定したコピーのコピーを前記複数の記憶装置間でランダムに分散するステップを含む

請求項24記載のプロセス。

【請求項36】 セグメントに対応する前記冗長情報は、2つ以上のセグメントに基づき、更に、

前記特定したセグメントに対応する前記冗長情報にしたがって、前記特定したセグメントを再構成するステップを含み、

前記ランダムに分散するステップは、前記再構成した特定セグメントを前記複数の記憶装置間で分散するステップを含む、

請求項34記載のプロセス。

【請求項37】 ビデオ・データ・ストリームを組み合わせる複合ビデオ・データを生成し、ビデオ・データを格納するための複数の記憶装置を備えた分散システムに前記複合ビデオ・データを格納するプロセスであって、前記記憶装置上に格納した前記ビデオ・データのセグメントおよび対応する冗長情報を、前記複数の記憶装置間でランダムに分散する際に、

前記複数の記憶装置から前記ビデオ・データ・ストリームを読み出すステップと、

前記ビデオ・データ・ストリームを組み合わせ、前記複合ビデオ・データを生成するステップと、

前記複合ビデオ・データをセグメントに分割するステップと、

前記複合ビデオ・データのセグメントおよび対応する冗長情報を前記複数の記憶装置間でランダムに分散するステップと、

からなるプロセス。

## 【発明の詳細な説明】

【0001】

## (発明の分野)

本発明は、マルチメディア・プログラムのキャプチャ、オーサリングおよび再生のためのコンピュータ・システム、ならびに分散型計算機システムに関するものである。

【0002】

## (背景)

コンピュータ・ネットワークを通じてデータの分散使用に対応するいくつかのコンピュータ・システム・アーキテクチャがある。これらのコンピュータ・アーキテクチャは、社内イントラネット、分散データベース・アプリケーション、およびビデオ・オン・デマンド・サービスのよう用途に用いられている。

【0003】

例えば、ビデオ・オン・デマンド・サービスは、典型的に、ユーザがムービー全体を要求しており、選択されたムービーはかなりの長さを有するという仮定の下で設計する。したがって、ビデオ・オン・デマンド・サーバは、数人の加入者による同一ムービーへの、恐らくは異なる時点における、読み取り専用アクセスに対応するように設計されている。このようなサーバは、通常、データを数個のセグメントに分割し、これらのセグメントを数台のコンピュータまたはコンピュータ・ディスクに順次分散する。この技法は一般にストライピング (striping) と呼ばれており、例えば、米国特許第5,473,362号、第5,583,868号、および第5,610,841号に記載されている。数台のディスクにムービー用データをストライプする際の1つの問題として、1台のディスクまたはサーバの故障のために、ムービー全てが失われる可能性があることがあげられる。何故なら、あらゆるムービーは少なくとも1つのセグメントが各ディスク上に書き込まれているからである。

【0004】

データ・ストレージにおける信頼性を高める一般的な技法に、ミラーリング (mirroring) と呼ばれるものがある。ミラーリングおよび順次ストライ

ピングを用いた混成システムが、米国特許第5, 559, 764号 (Chen et al. (チェンその他)) に示されている。ミラーリングは、各記憶装置のコピーを2つ維持すること、即ち、全てのデータについて一次ストレージおよび二次バックアップ・ストレージを有する必要がある。コピーは双方とも、負荷分散のために使用することができる。しかしながら、この技法を用いると、一次ストレージが故障した場合、その負荷全体が二次バックアップ・ストレージ上にかかることになる。

【0005】

数台のディスクにわたってデータを順次ストライプすることに伴う別の問題として、「コンボイ効果」(convoy effect) と呼ばれるものの確度が高くなることがあげられる。コンボイ効果が発生するのは、あるファイルからのデータ・セグメントに対する要求が1つのディスクに集中し、次いでディスクからディスクに循環する(「コンボイ」)する傾向があるからである。その結果、1つのディスクに特定の一度に要求の負担がかかり、一方他のディスクの負荷は軽くなる。また、ディスクに対する新しい要求は、そのいずれもがコンボイが処理されるのを待たなければならず、その結果、新たな要求に対するレイテンシが増大することになる。コンボイ効果を克服するためには、データをランダムにストライプすればよい。即ち、データ・ファイルのセグメントを、順次ではなく、ディスク間でランダムな順序で格納する。このようなシステムは、Proceedings of Multimedia' 96, pp144~150に、R. Tewari et al. (R. テワリその他) による "Design and Performance Tradeoffs in Clustered Video Servers" (クラスタ化ビデオ・サーバにおける設計および性能のトレードオフ) に記載されている。このようなシステムでもなお、1つのディスク上でのランダムな過剰負荷が生ずるが、これは一般的なデータ・アクセスのランダム性によるものである。

【0006】

これらのシステムには、多数の記憶装置および多数のアプリケーション間においてスケーラブルにかつ信頼性高く、多数の独立した高帯域データ・ストリーム

、特に、ビデオおよび関連するオーディオ・データのような等時性メディア・データを個別に転送可能なものはない。このようなデータ転送の必要性は、特に、マルチメディア・データのキャプチャ、オーサリング、および再生に対応するシステムでは困難である。特に、オーサリング・システムでは、データへのアクセスは、典型的に、より大きなデータ・ファイルの、クリップと呼ばれる、小断片を単位として行われる。これらのクリップは、どのようにデータを格納するかに関しては、任意に即ちランダムな順序でアクセスする場合が多く、効率的なデータ転送を行なうことが困難である。

(概要)

コンピュータ・ネットワークを用いて、多数のアプリケーションと接続してある多数の記憶装置上にデータをランダムに分散する。データをセグメントに分割する。各セグメントを記憶装置の1つに格納する。また、1つ以上のセグメントに基づく冗長情報も、基本となるセグメントとは異なる記憶装置上に格納する。冗長情報は、各セグメントのコピーとすればよく、あるいは2つ以上のセグメントに対して行なう排他的OR演算によって計算してもよい。セグメントまたは冗長情報を格納した各記憶装置の選択はランダムまたは疑似ランダムであり、他のデータ・セグメントが格納されている記憶装置とは独立することができる。冗長情報が2つ以上のセグメントに基づく場合、セグメントの各々を異なる記憶装置上に格納する。

【0007】

このデータ・セグメントのランダム分散によって、スケーラビリティ (scalability) および信頼性双方が向上する。例えば、セグメントにアクセスすることによってデータを処理するので、データ・フラグメント即ちクリップもデータの全てと同様に効率的に処理される。アプリケーションは、データ転送が効率的な場合にのみ、記憶装置からのデータ転送を要求することができ、更に記憶装置にリード要求を前処理するように要求することも可能である。コンピュータ・ネットワーク上での帯域幅利用は、クライアントおよび記憶装置間におけるデータ転送をスケジューリングすることによって最適化することができる。記憶装置の1つが故障した場合、その負荷も残りの記憶装置全体にランダムにそし

てほぼ均一に分散される。記憶装置の故障から復元する手順も備えることができる。

【0008】

記憶装置およびアプリケーションは、中央制御部を用いずに、独立して動作することも可能である。例えば、各クライアントは、ローカル情報のみを用いて、記憶装置との通信をスケジューリングすることができる。したがって、記憶装置およびアプリケーションのシステムに対する追加および除去が可能となる。その結果、システムは動作中でも拡張可能となる。

【0009】

冗長情報が1セグメントのコピーである場合、システム性能を向上させることができるが、記憶増大が犠牲となる。例えば、アプリケーションが選択したデータ・セグメントを要求する場合、要求は、要求のキューが最も短い記憶装置によって処理することができるので、多数の記憶装置上の多数のアプリケーションによって加えられる負荷におけるランダムな変動は統計的にそして一層等しく記憶装置全てにわたって均衡化される。

【0010】

この技法の組み合わせによって、多数の記憶装置および多数のアプリケーション間でスケーラブルにかつ信頼性高く多数の独立した高帯域幅データ・ストリームを転送することができるシステムが得られる。

【0011】

したがって、一態様では、分散データ記憶システムは、データを格納する複数の記憶装置を含み、記憶装置上に格納したデータのセグメントを、複数の記憶装置間で分散する。各セグメントに対応する冗長情報も、記憶装置間にランダムに分散する。

【0012】

冗長データが1セグメントのコピーである場合、各セグメントの各コピーは、記憶装置の異なるものに格納することができる。各セグメントの各コピーは、記憶装置の相対的指定 (relative specification) の関数として定義した確率分布にしたがって、複数の記憶装置の1つに割り当てること

ができる。分散データ記憶システムは、コンピュータ読み取り可能媒体を含むことができる。この媒体上にはコンピュータ読み取り可能なロジックを格納しており、データ・セグメントの指示を用いてコンピュータによってアクセス可能なセグメント・テーブルを定義し、セグメントのコピーを格納してある複数の記憶装置から、記憶装置の指示を検索する。複数の記憶装置は、コンピュータ・ネットワークに接続してある、第1、第2および第3記憶装置を含むことができる。

【0013】

別の態様では、コンピュータ用ファイル・システムは、コンピュータ上で実行するアプリケーションからの、要求に応答してコンピュータ・ネットワークを通じて独立したリモート記憶装置にアクセスし、前記記憶装置上に格納してあるデータを読み出すことを可能とする。前記データのセグメントおよび対応する冗長情報を前記複数の記憶装置間でランダムに分散する。冗長情報がセグメントのコピーである場合、ファイル・システムは、要求に応答してデータを読み出し、選択したデータの各セグメント毎に、セグメントを格納する記憶装置の1つを選択する。ファイル・システムは、他のセグメントおよび冗長情報から、失われたセグメントを再構成することができる。要求データの各セグメントは、当該セグメントのために選択した記憶装置から読み出される。選択した記憶装置からデータを受信すると、データをアプリケーションに供給する。このファイル・システムでは、複数の記憶装置上の要求の負荷が実質的に均衡化するように、記憶装置を選択することができる。セグメントのための記憶装置は、要求に対応する推定時間が最も短い、セグメントのための記憶装置の推定値に応じて選択することができる。

【0014】

更に特定すれば、ファイル・システムは、記憶装置の1つからデータを要求し、推定時間を示すことができる。第1記憶装置が要求を拒絶した場合、ファイル・システムは別の記憶装置からデータを要求することができ、別の推定時間を示すことができる。第2記憶装置が要求を拒絶した場合、ファイル・システムは第1記憶装置からデータを要求する。各記憶装置は、推定時間以内に当該記憶装置が要求に応じられないときに、データ要求を拒絶する。記憶装置は、推定時間以

内に当該記憶装置が要求に応じられるときに、要求を受け入れる。

【0015】

ファイル・システムは、選択した記憶装置からのデータ転送をスケジュールし、記憶装置が効率的にデータを転送するように、各セグメントを読み出すことができる。更に特定すれば、ファイル・システムは、選択した記憶装置からのデータ転送を要求し、待ち時間を示すことができる。選択した記憶装置がデータを転送する要求を拒絶した場合、他の記憶装置からデータを要求することができ、またはファイル・システムが後の時点で同じ記憶装置からデータを要求することも可能である。各記憶装置は、示された待ち時間以内に記憶装置からデータを転送することができない場合、データ転送の要求を拒絶する。記憶装置は、示された待ち時間以内に、選択した記憶装置がデータを転送することができる場合、データを転送する。

【0016】

別の態様では、コンピュータ用ファイル・システムは、コンピュータ上で実行するアプリケーションからの要求に応答して、コンピュータ・ネットワークを通じて独立したリモート記憶装置にアクセスし、前記記憶装置上にデータを格納することを可能とする。ファイル・システムは、データ格納要求に応答して、データを複数のセグメントに分割する。各セグメントは、1つ以上のセグメントに基づく冗長情報と共に、複数の記憶装置間でランダムに分散される。ファイル・システムは、アプリケーションに、データを格納するか否か確認する。

【0017】

このファイル・システムでは、冗長データがセグメントのコピーである場合、データのランダム分散は、各セグメント毎に、ランダムにかつ他のセグメントに選択した記憶装置には独立して、少なくとも2つの記憶装置を選択することによって行なうことができる。選択した記憶装置に、各セグメント毎にデータを格納するように要求することができる。ファイル・システムは、記憶装置の部分集合を選択することができ、選択した部分集合内にある記憶装置から、セグメントを格納する記憶装置を選択することができる。

【0018】

また、ファイル・システムの機能性は、他のアプリケーションによって、またはアプリケーション・プログラム・インターフェースを通じてアクセス可能なコード・ライブラリによっても得ることができる。したがって、別の態様は、記憶装置の選択およびネットワーク転送のスケジューリングを含む、リードまたはライト機能を実行するように実現したクライアントまたはプロセスである。別の態様は、記憶装置の選択およびネットワーク転送のスケジューリングを含むリードまたはライト機能を実行するために、これによって実現した記憶装置またはプロセスである。別の態様は、このような機能性を実現する分散コンピュータ・システムである。これらの動作は、クライアントまたは記憶装置によって、ローカル情報のみを用いて実行することができ、システムを容易に拡張可能とする。

【0019】

別の態様では、データを記憶するための複数の記憶装置を有する分散データ記憶システムにおいてデータを復元する。記憶装置の1つの故障が検出された場合、記憶装置上に格納したデータ・セグメントおよび冗長情報を複数の記憶装置間にランダムに分散する。データを復元するために、故障した記憶装置上にコピーを格納したセグメントを特定する。特定したセグメントに対応する冗長データを格納した記憶装置を特定する。冗長情報を用いて、特定したセグメントのコピーを再構成し、次いでこれらを複数の記憶装置間でランダムに分散する。このようなデータ復元は、ここに記載するファイル・システムまたは分散記憶システムのリードおよびライト機能性と組み合わせて用いることができる。

【0020】

別の態様では、ビデオ・データ・ストリームを組み合わせて複合ビデオ・データを生成し、ビデオ・データを格納するための複数の記憶装置を備えた分散システムに格納する。記憶装置上に格納したビデオ・データのセグメントのコピーを、複数の記憶装置間でランダムに分散する。ビデオ・データ・ストリームを含む記憶装置から読み出す。これらのビデオ・データ・ストリームを組み合わせて、複合ビデオ・データを生成する。複合ビデオ・データをセグメントに分割する。複合ビデオ・データのセグメントのコピーを、複数の記憶装置間でランダムに分散する。データの読み取りおよび格納は、ここに記載する技法を用いて行なう



ことができる。

(詳細な説明)

添付図面に関連付けて読むべき以下の詳細な説明では、本発明の実施形態の例を明示する。ここに引用する全ての引例は、ここで言及することによって明示的に本願にも含まれるものとする。

#### 【0021】

多数のアプリケーションおよび多数の記憶装置間において、データ転送、特に、モーション・ビデオおよび付随するオーディオ、ならびにその他の時間的に連続するメディアのような、多数の独立した高帯域幅時間厳守データ・ストリームの転送に対応するスケラブルで信頼性の高い分散システムの設計において、いくつかの問題点が生ずる。このようなシステムでは、例えば、モーション・ビデオ・プログラムのオーサリングを行なうために用いるアプリケーションは、数個の記憶装置にわたって分散している可能性がある数個の異なるファイルの数個の小部分にランダムにアクセスする場合がある。数個のアプリケーションが同じデータに即時かつ同時のアクセスを要求する場合もあり、いずれのアプリケーションもいずれの時点でもいずれのメディア・ピースにもアクセスすることができなければならない。放送またはその他時間厳守再生に用いるシステムでは、フォールト・トレランスも望ましい。最後に、システムは、新たな記憶装置や新たなアプリケーションの追加を、システムの動作中においても簡単にできるように、拡張可能かつスケラブルでなければならない。このようなシステムに望ましいその他の特性には、故障までの平均時間が長いこと、単一点で故障しないこと、迅速にかつ動作しながらでも修理可能であること、動作を中断することなく記憶装置の故障に対する耐性があること、そして失われたデータを復元可能であることが含まれる。

#### 【0022】

一実施形態において、本システムは、コンピュータ・ネットワークによって多数の別個で独立したデータ記憶用記憶装置に接続した多数のアプリケーションを含む。データはセグメントに分割する。各セグメント毎の冗長情報を決定し、セグメントおよびその冗長情報を別々の記憶装置上に格納する。セグメントに対す

る記憶装置の選択はランダムまたは疑似ランダムであり、直前のセグメントのよ  
うな、他のセグメントに選択した記憶装置とは独立とすることもできる。冗長情  
報およびランダムなデータ分散双方によって、アプリケーションおよびストレ  
ージ間において双方向にデータを効率的に転送するシステムの機能が向上し、フォ  
ールト・トレランス性が高くなる。

【0023】

冗長情報は、セグメントのコピーとすることができる。このセグメントの複製  
により、システムは、要求のキューが最も短い記憶装置を選択することによる等  
のように、特定のアプリケーションがどの記憶装置をアクセスするのも制御す  
ることが可能となる。その結果、負荷のランダムな変動は、記憶装置全てにほぼ  
均等に分散される。

【0024】

また、アプリケーションは、転送が効率的な場合にのみ、記憶装置にデータ転  
送を要求することも可能である。ネットワーク上での通信を適切にスケジュー  
ルすることにより、ネットワークの輻輳を低減することができ、ネットワークの帯  
域幅を一掃効率的に使用することができる。各クライアントにローカル情報を使  
用して記憶装置との通信をスケジュールさせることによって、中央制御点を不要  
とすることも可能である。

【0025】

図1Aは、一例としてのコンピュータ・システム40を示す。このコンピュ  
ータ・システムは複数の記憶装置42を含む。記憶装置とは、ディスクのような不  
揮発性コンピュータ読み取り可能媒体を備え、その上にデータを格納することが  
できるデバイスのことである。また、記憶装置は、高速な、典型的に不揮発性の  
メモリも有し、その中に媒体からのデータを読み出す。各記憶装置は、それ自体  
の独立したコントローラも有し、リードおよびライト・アクセスを含みこれらに  
は限定されない、媒体上に格納してあるデータへのアクセス要求に応答する。例  
えば、記憶装置42は、サーバ・コンピュータとすることができ、サーバのファ  
イル・システム内にあるデータ・ファイルにデータを格納する。コンピュータ・  
システム40には任意の数の記憶装置があり得る。

## 【0026】

アプリケーション44は、コンピュータ・ネットワーク46を通じて記憶装置への要求を介して、記憶装置42へのアクセスを要求するシステムである。記憶装置42は、コンピュータ・ネットワーク46を通じて、アプリケーション44にデータを送出したり、あるいはアプリケーション44からのデータを受け取る。アプリケーション44は、ディジタルまたはアナログ・ソースから受け取ったデータを捕獲し、記憶装置42上にデータを格納するシステムを含むことができる。また、アプリケーション44は、マルチメディア・プログラムのオーサリング、処理または再生用システムのように、記憶装置からデータを読み出すシステムも含むことができる。他のアプリケーション44は、種々の故障回復タスクを実行することができる。また、アプリケーション44を「クライアント」と呼ぶ場合もある。1つ以上のカタログ・マネージャ49も使用可能である。カタログ・マネージャとは、アプリケーション44によるアクセスが可能なデータベースであり、記憶装置42上で利用可能なデータに関する情報を保持している。この実施形態は、1997年10月23日付けのPCT公開WO97/39411に示されるような放送ニュース・システムを実現するために用いることも可能である。

## 【0027】

記憶装置42に格納するデータは、セグメントに分割する。1つ以上のセグメントに基づいて、冗長情報を作成する。例えば、各セグメントをコピーすればよい。その結果、少なくとも2つの記憶装置42上に格納各セグメントを格納することになる。あるいは、2つ以上のセグメントの排他的ORによって、冗長情報を作成することも可能である。各セグメントは、その冗長情報とは異なる記憶装置42に格納する。セグメントおよびその冗長情報を格納する記憶装置の選択はランダムまたは疑似ランダムであり、他のデータ・セグメントを格納する記憶装置からは独立とすることも可能である。一実施形態では、同じ記憶装置には、2つの連続するセグメントを格納しない。セグメントおよびその冗長情報を格納する記憶装置を選択するための確率分布は、記憶装置全体にわたって均一とするとよく、容量、帯域幅、およびレイテンシというような、記憶装置の仕様は同様と

する。この確率分布は、各記憶装置の仕様の関数とすることも可能である。データ・セグメントおよび対応する冗長情報のランダム分散により、スケーラビリティおよび信頼性双方が向上する。

【0028】

データ・セグメントのコピーのランダム分散の一例を図1Aに示す。図1Aにおいて、w, x, yおよびzで示す4つの記憶装置42が、1, 2, 3および4で示す4つのセグメントに分割したデータを格納する。セグメントおよびそのコピーのランダム分散の一例を示す。ここで、セグメント1および3を記憶装置wに格納し、セグメント3および2を記憶装置xに格納し、セグメント4および1を記憶装置yに格納し、セグメント2および4を記憶装置zに格納する。

【0029】

図1Bは、セグメントおよびその対応する冗長情報を記憶装置間でランダムに分散した一実施形態を示す。図1Bでは、w, x, yおよびzで示す4つの記憶装置42が、1, 2, 3および4で示す4つのセグメントに分割したデータを格納する。セグメントに対する冗長情報は、1つ以上のセグメントに基づくことができる。この例では、ここでは「冗長セット」と呼ぶものにおいて、2つのセグメントを用いる。冗長セットにおけるセグメントi, jの排他的ORを計算して、冗長情報 $R_{ij}$ を得る。冗長情報 $R_{ij}$ およびセグメントiの排他的ORにより、セグメントjを生成する。同様に、冗長情報 $R_{ij}$ およびセグメントjの排他的ORにより、セグメントiを生成する。冗長セット内の各セグメントおよび冗長情報は、異なる記憶装置上に格納する。セグメントおよび冗長情報のランダム分散の一例を図1Bに示す。ここでは、セグメント3および4に対する冗長情報 $R_{3,4}$ を記憶装置w上に格納し、セグメント2および3を記憶装置x上に格納し、セグメント1を記憶装置y上に格納し、セグメント4および冗長情報 $R_{1,2}$ を記憶装置z上に格納する。また、フォールト・トレランスの技術分野においては公知である他の多くの技法を用いて、冗長情報を作成することも可能である。

【0030】

冗長情報がセグメントのコピーである場合、セグメントのランダム分散は、セグメント・テーブル90、または図2Aに示すようなカタログに内において表わ

## 【0032】

各セグメント・テーブルまたはファイル・マップは、他のセグメント・テーブルとは別個に格納することもできる。セグメント・テーブルは、カタログとして、一緒に格納することができる。カタログは、個々のクライアントにおいて、中央データベースにあるカタログ・マネージャ49上に格納することができ、あるいはいくつかのデータベースまたはクライアント間で分散することも可能である。例えば、異なる種類のメディア・プログラム毎に、別個のカタログを維持することも可能である。例えば、放送ニュースの編成は、スポーツ・ニュース、天気予報、ヘッドライン・ニュース等に対して別個のカタログを有することができる。また、カタログは、他のデータと同様に記憶装置上に格納することも可能である。例えば、各クライアントは、乱数発生器のシードを用いて、カタログにアクセスすることができる。このようなカタログは、例えば、ネットワーク・ブロードキャスト・メッセージを全てのカタログ・マネージャまたはクライアントに送り、個々のセグメント・テーブルのカタログのコピーを得ることによって、他のクライアントが特定してデータにアクセスしたり、または復元要求を処理することができる。

## 【0033】

データ・セグメントにアクセスするためには、各セグメントは一意の識別子を有していなければならない。セグメントのコピーは同じ一意の識別子を有することができる。2つ以上のセグメントに基づく冗長情報はそれ自体の識別子を有する。セグメントに対する一意の識別子は、ファイルのようなソースに対する一意の識別子、およびセグメント番号の組み合わせである。ソースまたはファイルに対する一意の識別子は、例えば、システム時間、またはデータをソースから捕獲した時に判定したその他の一意の識別子によって、あるいはファイルの作成時に決定することができる。以下で説明するが、ファイル・システムは、カタログ・マネージャにアクセスし、各ソース毎のセグメント・テーブル、またはセグメント識別子ならびにセグメントおよび冗長情報を格納してある記憶装置をリストに纏めたファイルを得ることができる。また、各記憶装置は、別個のファイル・システムも有することができ、これには、セグメント識別子、およびこれらを格納

してある記憶装置上の位置のディレクトリを収容する。クライアントが実行するアプリケーション・プログラムは、ソースまたはファイルの識別子、および恐らくは当該ソースまたはファイル内のある範囲のバイトを用いて、クライアントのファイル・システムからデータを要求することができる。すると、クライアントのファイル・システムは、一意のセグメント識別子を用いて、そのソースまたはファイルのセグメント・テーブルを突き止め、どのセグメントにアクセスする必要があるかについて判定を行い、各セグメント毎にデータを読み出す記憶装置を選択する。

【0034】

再度図1Aおよび図1Bを参照する。アプリケーション44が記憶装置42の1つにおいて選択したデータ・セグメントへのアクセスを要求する場合、記憶装置はその要求をキュー48上に置く。キュー48は、記憶装置が維持している。アプリケーションはこのような要求を、互いにまたはあらゆる中央制御部にも独立して行なうことができ、このためシステムは容易にスケーラブルとなっている。冗長情報がセグメントのコピーである場合、要求を送る先の記憶装置の選択は、多数の記憶装置42上の多数のアプリケーション44によって加えられる負荷のランダムな変動が、記憶装置42全てにわたって統計的にかつ一層等しく均衡化されるように、制御することができる。例えば、アプリケーション44からの各要求は、最も短い要求のキューを有する記憶装置によって処理することができる。あらゆる種類の冗長情報でも、アプリケーションおよび記憶装置間のデータの転送は、ネットワークの輻輳を低減するようにスケジューリングすることができる。データ要求は、2段階で実行することができる。即ち、データをディスクから記憶装置上のバッファに転送する予備リード要求と、ネットワークを通じてバッファからアプリケーションにデータを転送するネットワーク転送要求である。これら2つの異なる要求を処理するために、キュー48は、ディスク・キューおよびネットワーク・キューを含むとよい。

【0035】

このランダムに分散したデータ・セグメントおよび対応する冗長情報の組み合わせ、ならびにネットワークを通じたデータ転送のスケジューリングにより、多

数の記憶装置および多数のアプリケーション間でスケーラブルにかつ信頼性高く双方向に、多数の独立した高帯域幅データ・ストリームを転送することができるシステムが得られる。セグメントのコピーを冗長データとして用いることにより、リード・アクセスのための記憶装置の選択は、記憶装置の相対的な負荷に基づいて行なうことができ、処理能力の向上を図ることができる。

【0036】

次に図3を参照し、いくつかの記憶装置にランダムに分散してデータ・セグメントの多数のコピーを格納するプロセス例について、更に詳しく説明する。1つ以上のセグメントに基づく冗長情報を用いたプロセス例については、図24に関連付けて以下で説明する。以下の説明は、モーション・ビデオ・データのリアル・タイム・キャプチャに基づくものである。この例は、オーディオのような他の時間的に連続するメディア、静止画像やテキストのような離散メディア、または知覚データのようなその他のデータをも含むがこれらには限定されない、他の形態のデータに一般化することができる。

【0037】

米国特許第5,640,601号および第5,577,190に記載されているように、どのようにしてリアル・タイム・モーション・ビデオ情報をコンピュータ・データ・ファイルに取り込むかは、一般的に周知である。この手順は、取り込んだデータをセグメントに分割し、コピーし、セグメントのコピーを記憶装置間でランダムに分散するためのステップを含むように変更することも可能である。最初に、ステップ120において、キャプチャ・システムはセグメント・テーブル90A（図2A）を作成する。典型的に、取り込むデータ・ストリームへのオフセットに、各画像をマップする画像インデックスも作成する。インデックス化したデータは、例えば、フィールドまたはフレームに対応することができる。インデックスは、オーディオのような別の種類のデータについては、時間期間のような他のサンプル境界を引用することができる。また、キャプチャ・システムは、利用可能な記憶装置のリストも取得する。どの記憶装置が利用可能かを識別する方法については、図10ないし図12に関連付けて以下で更に詳しく説明する。

## 【0038】

ステップ121において、キャプチャ・システムはデータ・セグメントを作成する。セグメントのサイズは、例えば、モーション・ビデオ情報では1/4メガバイト、1/2メガバイト、または1メガバイトとするとよい。オーディオ情報は、例えば、1/4メガバイトのようなサイズを有するセグメントに分割するとよい。セグメント・サイズの格納および送信の分割に対する整合を得るためには、可能であれば、セグメントのサイズは、未圧縮または固定データ・レート、ディスク・ブロックおよびトラック・サイズ、メモリ・バッファ・サイズ、ネットワーク・パケット（例えば、64K）および/またはセル・サイズ（例えば、ATMでは53バイト）の倍数に関係付けるとよい。データが未圧縮かまたは固定レートの圧縮を用いて圧縮されている場合、セグメントは、時間的なサンプル境界で分割すれば、画像インデックスとセグメント・テーブルとの間に整合を得ることができる。概して言えば、セグメント・サイズは、システム・オーバーヘッドを低減するためには、大きくする方がよい。システム・オーバーヘッドは、セグメントが小さい程増大する。一方、記憶装置の全てにデータが分散されないような格納データ量およびセグメント・サイズである場合、コンボイ効果が発生する確率は高くなる。加えて、セグメント・サイズが大きい程、ディスク要求およびネットワーク要求を完了するためのレイテンシも増大する。

## 【0039】

次に、ステップ122において、キャプチャ・システムは、選択したセグメントを格納するために利用可能な記憶装置のリストから、少なくとも2つの記憶装置42を選択する。一方のセグメントのコピーのための記憶装置の選択は、ランダムまたは疑似ランダムである。この選択は、直前または直後のセグメントに対して行なう選択とは独立とすることも可能である。また、選択を行なった記憶装置の集合は、利用可能な記憶装置全ての部分集合とすることも可能である。1組の記憶装置の選択は、各ソースまたはファイル毎にランダムまたは疑似ランダムとすることもできる。この部分集合のサイズは、各記憶装置が少なくとも1つの異なるデータ・セグメントを有し、コンボイ効果の発生確度を極力低下させるようにしなければならない。即ち、データは、1組の記憶装置の数の少なくとも2



倍の長さ（セグメントにおいて）としなければならない。部分集合のサイズも、いずれの所与の時点においても、部分集合内の2つ以上の記憶装置が故障する、即ち、二重故障が発生する確度を低下させるためには、制限しなければならない。例えば、5つの内2つの記憶装置が故障し得る確率は、100個の内2つの記憶装置が故障し得る確率よりも低く、したがって、データを分散する記憶装置の数を制限しなければならない。しかしながら、性能と部分集合のサイズとの間にはトレードオフがある。例えば、100個の記憶装置から10個の部分集合をランダムに選択して用いる場合、100個の記憶装置の内2つが故障すると、ファイルの10パーセントが悪影響を受ける。部分集合がないと、典型的に、ファイルの100パーセントが悪影響を受けることになる。

#### 【0040】

二重故障、即ち、2つ以上の記憶装置が故障するという稀な確度では、データ・セグメントが失われる可能性がある。標準的なビデオ・ストリームでは、1セグメントの損失のために、プログラム素材1分において1または2フレームが失われる可能性がある。所与のソースまたはファイルに対するこのような故障の頻度は、その帯域幅および記憶装置数の関数である。即ち、

$s$  = メガバイト (MB) 単位での失われたデータのサイズ

$n$  = 記憶装置の初期数

$b$  = MB 単位での1秒当たりの記憶装置の平均帯域幅

$MTBF$  = 故障間平均時間

$MTTR$  = 故障または交換までの平均時間

$MTDF$  = 二重故障不良の平均時間

$SMTBF$  = 故障間の全システム平均時間

とすると、

#### 【0041】

【数1】

$$SMTBF = \frac{MTBF}{n} \quad \text{および} \quad MTDF = \frac{1}{MTTR} * \frac{MTBF}{n} * \frac{MTBF}{(n-1)}$$

となる。

一例として、100個の記憶装置を有し、各々50ギガバイトの容量を有するシステムでは、MTTRが1時間、MTBFが1000時間即ち6週間の場合、115年で二重故障不良が発生する確度となる。MTTRが24時間に増大すると、4.8年で二重故障不良が発生する確度となる。

#### 【0042】

再び図3を参照する。2つの記憶装置を選択した後、ステップ124において現セグメントを選択した記憶装置の各々に送り格納する。これらのライト要求は、順次ではなく、非同期とすることも可能である。次に、キャプチャ・システムは、ステップ126において、記憶装置がセグメントの格納完了を承認するのを待つ。捕獲しながらデータをリアル・タイムで格納する場合、ステップ124におけるデータ転送は、以下で更に詳しく論ずるリード動作と同様、2段階で行なうとよい。即ち、最初にクライアントが記憶装置に、データを格納するための空きバッファを準備するように要求することができる。記憶装置は、バッファが利用可能な推定時間を回答することができる。この推定時間に到達した場合、キャプチャ・システムは、記憶装置にデータを受け取るように要求することができる。次いで、記憶装置はそのバッファ内のデータを受け取り、そのバッファ内のデータをその記憶媒体に転送し、承認をキャプチャ・システムに送ればよい。

#### 【0043】

キャプチャ・システムが承認を受け取る前に、時間切れとなった場合、セグメントを再度同じ記憶装置または異なる記憶装置に送ることもできる。他のエラーも、キャプチャ・システムによって処理することができる。選択した記憶装置上での格納成功を確実にする動作を、セグメントの各コピー毎に別個のスレッドによって行なうことも可能である。

#### 【0044】

データを記憶装置に格納することに成功した後、ステップ127においてキャプチャ・システムはセグメント・テーブル90を更新する。キャプチャが完了したとステップ128において判断した場合、プロセスは終了する。それ以外の場合、ステップ121に戻って、次のセグメントのためにプロセスを繰り返す。セ

グメント・テーブルは、例えば、キャプチャ・システムの主メモリ内に、ファイル・システムの一部として維持することができる。本例では、キャプチャ・システムはセグメント・テーブルおよび記憶装置の選択を管理するが、カタログ・マネージャ49のようなシステムの他の部分が、これらのアクティビティを調整することも同様に可能である。ステップ129において、更新したセグメント・テーブルは、例えば、カタログ・マネージャに送ることができる。あるいは、カタログ・マネージャは、蓄積したシステム動作の知識を用いることによって、セグメント・テーブルを生成することもでき、更に要求に応じてこのテーブルをキャプチャ・システムに送ることができる。

【0045】

図4は、記憶装置がどのようにして捕獲したデータ・セグメントまたは冗長情報を格納するのかについて更に詳細に記載したフローチャートである。ステップ140において、記憶装置は、キャプチャ・システムからデータ・セグメントを受け取り、記憶装置のバッファにデータを格納する。記憶装置は格納のためにデータ・ファイルを用いると仮定すると、ステップ142において記憶装置はデータ・ファイルを開き、ステップ144においてデータ・ファイルにデータを格納する。カタログ・マネージャは、セグメントを格納するロケーションを指定することも可能である。データは、既存のデータ・ファイルに添付することができ、あるいは別個のデータ・ファイルに格納することもできる。先に論じたように、ステップ145において、記憶装置またはカタログ・マネージャは、各セグメント毎の一意の識別子を用いることにより、更にセグメント識別子を記憶装置上のそのロケーションにマップするテーブルを格納することにより、セグメントを追跡することも可能である。このテーブルは、記憶装置上におけるデータ・ファイルの抜き取り (a b s t r a c t i o n) を実現することができる。記憶装置がいつデータを実際にその主ストレージに書き込むかは、他のアプリケーションで保留となっているその他のリードおよびライト要求に左右される。これら同時要求の管理については、以下で更に詳しく扱うことにする。次に、ステップ146において、ファイルを閉じることができる。ステップ148において、承認をキャプチャ・システムに送ることができる。

## 【0046】

図3および図4のプロセスが完了したなら、捕獲したデータを、各セグメント毎に少なくとも2つのコピーを用いて、数個の記憶装置にランダムに分散する。多数のアプリケーションがこのデータへのアクセスを要求することができる。このアクセスを行なう態様は、恐らくランダムである。したがって、いずれの記憶装置も、多数のアプリケーションから、記憶装置上に格納してあるファイルからデータを読み出す要求およびファイルにデータを書き込む要求を多数受ける可能性があることは明白である。要求を管理するために、前述のように、要求キュー48を記憶装置42の各々によって維持する。実施形態の一例についての以下の説明では、記憶装置は2つのキュー、即ち、一方はディスク・アクセス要求のため、他方はネットワーク転送要求のために維持する。これらのディスクおよびネットワーク・キューの一実施形態については、以下で図19に関連付けて更に詳しく説明する。

## 【0047】

クライアント44上で実行するアプリケーション・プログラムがデータを要求した場合、各データ・セグメントを少なくとも2つの記憶装置上に格納する際に、要求を満たすように記憶装置を選択する。要求データに対するセグメント・テーブル90をこの目的のために用いる。記憶装置の選択は、データを要求するアプリケーション・プログラムによって、アプリケーション・プログラムを実行するクライアントのファイル・システムによって、記憶装置間の調整によって、またはカタログ・マネージャのような別のアプリケーションによって行なうことができる。選択はランダムでも疑似ランダムでもよく、あるいは最低使用頻度アルゴリズムに基づいて、または記憶装置のキューの相対的な長さに基づいて行なうことも可能である。利用可能な記憶装置上のキューの相対的な長さに基づいて記憶装置を選択することによって、多数のアプリケーションの負荷を記憶装置の集合全体に一層等しく分散することが可能となる。このような選択については、図16ないし図18に関連付けて以下で更に詳しく説明する。

## 【0048】

これより、特定のな一実施形態の詳細について更に説明する。この目的のため

には、記憶装置42は、サーバまたは独立制御記憶装置として実現するとよく、アプリケーション44をクライアントと呼ぶことにする。クライアントは種々のタスクを実行するアプリケーション・プログラムを実行することができる。サーバまたはクライアントを実現するのに適したコンピュータ・システムは、典型的に、主ユニットを含み、これは通常バスまたはスイッチのような相互接続機構を介してメモリ・システムに接続するプロセッサを含む。サーバおよびクライアント双方は、これらをコンピュータ・ネットワークに接続するネットワーク・インターフェースも有する。ネットワーク・インターフェースは、フォールト・トレランスに対応するために冗長化するとよい。また、クライアントは、ディスプレイのような出力デバイスや、キーボードのような入力デバイスを有することができる。入力デバイスおよび出力デバイスは双方とも、相互接続機構を介して、プロセッサおよびメモリ・システムに接続することができる。

#### 【0049】

尚、1つ以上の出力デバイスをクライアント・システムに接続してもよいことは理解されよう。出力デバイスの例には、陰極線管（CRT）ディスプレイ、液晶ディスプレイ（LCD）、プリンタ、モデムまたはネットワーク・インターフェースのような通信デバイス、ならびにビデオおよびオーディオ出力が含まれる。また、1つ以上の入力デバイスをクライアント・システムに接続してもよいことも理解されよう。入力デバイスの例には、キーボード、キーパッド、トラックボール、マウス、ペンおよびタブレット、モデムまたはネットワーク・インターフェースのような通信デバイス、ビデオおよびオーディオ・ディジタイザ、ならびにスキャナが含まれる。尚、本発明は、コンピュータ・システムと組み合わせて用いる特定の入力デバイスにも出力デバイスにも、更にはここに記載されるものにも限定されないことは理解されよう。

#### 【0050】

コンピュータ・システムは、「C」および「C++」プログラミング言語のような、高級コンピュータ・プログラミング言語を用いてプログラム可能な汎用コンピュータ・システムとすることができる。また、コンピュータ・システムは、特別にプログラムした特殊目的ハードウェアとすることもできる。汎用コンピュ

ータ・システムでは、プロセッサは典型的に市販のプロセッサであり、その例には、Intel（インテル社）から入手可能な、MMX技術を用いたPentium IIのようなx86シリーズ・プロセッサや、AMDやCyrixから入手可能な同様のデバイス、Motorolaから入手可能な680X0シリーズ・マイクロプロセッサ、Digital Equipment Corporation（ディジタル・エキップメント社）から入手可能なAlphaシリーズ、およびIBMから入手可能なPowerPCプロセッサがある。多くのその他のプロセッサも利用可能である。このようなマイクロプロセッサは、オペレーティング・システムと呼ばれるプログラムを実行することができ、その例には、Windows NT、Windows 95、UNIX、IRIX、Solaris、DOS、VMS、VxWorks、OS/Warp、Mac OS System 7およびOS 8オペレーティング・システムがある。オペレーティング・システムは、他のコンピュータ・プログラムの実行を制御し、スケジューリング、デバッグ、入出力制御、コンパイル、ストレージ割り当て、データ管理およびメモリ管理、ならびに通信制御および関連サービスを行なう。プロセッサおよびオペレーティング・システムは、高級プログラミング言語でアプリケーション・プログラムを各ためのコンピュータ・プラットフォームを定義する。

【0051】

各サーバは、大量の主メモリ、例えば、32メガバイトよりも遥かに多いメモリや、例えば、数ギガバイトのディスク容量を有する安価なコンピュータを用いて実現することができる。ディスクは、1つ以上の単体ディスク、または独立したディスクの冗長アレイ（RAID）、あるいはその組み合わせとすることができる。例えば、サーバは、Windows NTまたはVxWorksのようなリアル・タイム・オペレーティング・システムを有する、Pentiumまたは486マイクロプロセッサを用いたシステムとすることができる。オーサリング・システム、キャプチャ・システム、および再生システムは、この主の製品のために当技術分野において現在用いられているプラットフォームを用いて実現することができる。例えば、マサチューセッツ州TewksburyのAvid Technology, Inc.（アビッド・テクノロジー社）からのMEDIACO

MPOSERは、PowerPCマイクロプロセッサおよびMac OS System 7オペレーティング・システムを有するApple Computer, Inc. (アップル・コンピュータ社)からのPower Macintoshコンピュータを用いている。Windows NTオペレーティング・システムを備え、IntelからのMMX技術を用いたPentium IIプロセッサに基づくシステムも使用可能である。再生システムの例には、コロラド州BoulderのPluto Technologies International Inc. (プルート・テクノロジーズ・インターナショナル社)からの「SPACE」システム、またはMachintoshプラットフォームを用いたAvi TechnologyからのAIRPLAYシステムが含まれる。カタログ・マネージャは、Informixデータベースのような、適当なデータベース・システムに対応するいずれかのプラットフォームを用いて実現することができる。同様に、システム内で利用可能なデータの種類の追跡するアセット・マネージャも、このようなデータベースを用いて実現することができる。

#### 【0052】

コンピュータにおけるメモリ・システムは、典型的に、コンピュータ読み取り可能および書き込み可能不揮発性記録媒体を含み、その例には、磁気ディスク、光ディスク、フラッシュ・メモリおよびテープがある。ディスクは、フロッピー・ディスクまたはCD-ROMのようなりムーバブルでも、ハード・ドライブのように固定でもよい。ディスクは多数のトラックを有し、その中に、典型的に二進形態、即ち、一連の1および0として解釈した形態で、信号を格納する。このような信号は、マイクロプロセッサが実行するアプリケーション・プログラム、またはディスク上に格納してありアプリケーション・プログラムが処理する情報を定義することができる。典型的に、動作において、プロセッサは不揮発性記録媒体から集積回路メモリ素子にデータを読み出す。集積回路メモリ素子は、典型的に、ダイナミック・ランダム・アクセス・メモリ(DRAM)またはスタティック・メモリ(SRAM)のような、揮発性ランダム・アクセス・メモリである。集積回路メモリ素子は、ディスクよりも、プロセッサによる情報への速いアクセスを可能とする。プロセッサは通常集積回路メモリ内でデータを操作し、次いで

処理が完了すると、データをディスクにコピーする。ディスクと集積回路メモリ素子との間のデータ移動を管理する機構には、様々なものが公知であり、本発明はそれに限定されるものではない。また、本発明は特定のメモリ・システムにも限定される訳ではないことも理解されよう。

【0053】

本発明は、特定のコンピュータ・プラットフォーム、特定のプロセッサ、特定の高級プログラム言語にも限定される訳ではないことは理解されよう。加えて、コンピュータ・システムは、マイクロプロセッサ・コンピュータ・システムとしてもよく、あるいはコンピュータ・ネットワークを通じて接続した多数のコンピュータを含むことも可能である。

【0054】

前述のように、各記憶装置42は、サーバを通じてアクセスされた場合、各アプリケーション44はファイル・システムを有することができる。ファイル・システムは、典型的に、オペレーティング・システムの一部であり、データのファイルを維持する。ファイルは、有名論理構造であり、ファイル・システムによって定義されかつ実現され、データの論理レコードの名称およびシーケンスを、物理的なストレージ・メディア上のロケーションにマップする。ファイル・システムは、データの物理的ロケーションをアプリケーション・プログラムからマスクするが、ファイル・システムは通常隣接するブロック内にある1つのファイルのデータを物理的ストレージ・メディアに格納しようとする。ファイルは、種々のレコード型に特定して対応することができ、あるいは未定義として、アプリケーション・プログラムに解釈または制御させることも可能である。ファイルは、その名称またはその他の識別子でアプリケーション・プログラムによって呼ばれ、オペレーティング・システムが定義するコマンドを用いて、ファイル・システムを通じてアクセスされる。オペレーティング・システムは、ファイルを作成し、ファイルを開き、ファイルに書き込み、ファイルを読み取り、ファイルを閉じるための基本的ファイル動作を行なう。これらの動作は、ファイル・システムに応じて、同期でも非同期でもよい。

【0055】



ここで記載する場合、ファイルまたはソースのデータはセグメント単位で格納し、そのコピーまたはその他の形態の冗長情報を多数の記憶装置の間でランダムに分散する。

【0056】

殆どのファイル・システムについて概して言えば、ファイルを作成するためには、オペレーティング・システムは最初に、ファイル・システムが制御するストレージ内において、空間を特定する。次に、新たなファイルに対するエントリをカタログ内に作成する。カタログは、利用可能なファイルの名称およびファイル・システム内におけるそれらのロケーションを示すエントリを含む。ファイルの作成は、一定の利用可能な空間をファイルに割り当てることを含むこともある。一実施形態では、ファイルに対してセグメント・テーブルを作成することができる。ファイルを開くと、典型的に、ハンドルがアプリケーション・プログラムに戻され、これを用いてそのファイルにアクセスする。ファイルを閉じると、ハンドルは無効となる。ファイル・システムはハンドルを用いて、ファイルに対するセグメント・テーブルを識別する。

【0057】

ファイルにデータを書き込むために、アプリケーション・プログラムはコマンドをオペレーティング・システムに発行し、オペレーティング・システムは、ファイル名、ハンドルまたはその他の記述子のようなファイルのインディケータ、およびファイルに書き込む情報の双方を指定する。概して言えば、ファイルのインディケータが与えられると、オペレーティング・システムはディレクトリを探索して、ファイルのロケーションを見つけ出す。データは、ファイル内の既知のロケーションまたはファイルの終端に書き込むことができる。ディレクトリ・エントリは、ライト・ポインタと呼ばれる、ポインタを現在のファイル終端に格納することができる。このポインタを用いて、ストレージの次に利用可能なブロックの物理的ロケーションを計算し、そのブロックに情報を書き込むことができる。ライト・ポインタをディレクトリ内において更新し、新たなファイル終端を示すことができる。一実施形態では、ライト動作は、ファイルのセグメントのコピーを記憶装置間でランダムに分散し、ファイルに対するセグメント・テーブルを更

新する。また、ライト動作も、セグメントおよび対応する冗長情報を異なる記憶装置に格納することができる。

【0058】

ファイルからデータを読み出すためには、アプリケーション・プログラムは、オペレーティング・システムにコマンドを発行し、ファイルのインディケータ、およびアプリケーションに割り当てられているメモリ・ロケーションの内リード・データが位置するメモリ・ロケーションを指定する。概して言えば、ファイルのインディケータが与えられると、オペレーティング・システムはそのディレクトリを探索して関連するエントリを求める。アプリケーション・プログラムは、用いるファイルの開始からのオフセットを指定することができ、あるいはシーケンシャル・ファイル・システムでは、ディレクトリは次に読み取るデータ・ブロックへのポインタを与えることができる。一実施形態では、記憶装置の選択やデータ転送のスケジューリングは、クライアントのファイル・システムのリード動作の一部として実現する。

【0059】

クライアントは、既定のアプリケーション・プログラミング・インターフェース（API）と共にファイル・システムまたは特殊コード・ライブラリを用いて、ファイルの一部に対する要求を、選択した記憶装置からのデータ・セグメントに対する要求に変換する。記憶装置は、クライアントのファイル・システムとは完全に分離した、それ自体のファイル・システムを有することも可能である。記憶装置上のセグメントは全て、例えば、記憶装置の単一ファイル内に格納することができる。あるいは、クライアント・ファイル・システムは、ネットワーク上で記憶装置を生ストレージ（raw storage）として用い、カタログ・マネージャおよびセグメント・テーブルを用いてファイルの抜き取りを実現することも可能である。ファイルのセグメント・テーブルは、各セグメントに選択した記憶装置上における当該セグメントのロケーションも示すことができる。

【0060】

ファイル・システムを用いる主要な利点の1つは、アプリケーション・プログラムにとっては、ファイルは、物理的な記憶媒体やオペレーティング・システム

がデータを格納するために用いる当該媒体上のロケーションについて何の懸念もなく、作成し、開き、書き込み、読み出し、そして閉じることができる論理構造であるという点にある。ネットワーク・ファイル・システムでは、ファイル・システムは、種々の記憶装置からの指定ファイルからのデータ要求を管理し、アプリケーション・プログラムは、データが格納されている物理的ストレージまたはコンピュータ・ネットワークに関する詳細を知る必要は全くない。記憶装置がそれ自体の独立したファイル・システムを有する場合、クライアント・ファイル・システムも、記憶装置の記憶機構の詳細を知る必要はない。記憶装置は、例えば、Windows NTファイル・システムと関連するファイル・システム、またはVxWorksのようなリアル・タイム・オペレーティング・システムのファイル・システム、または非同期動作を許すファイル・システムを用いることができる。

#### 【0061】

記憶装置は、クライアント、およびオプションとしてコンピュータ・ネットワークを用いてカタログ・マネージャと相互接続されている。コンピュータ・ネットワークは、1組の通信チャネルであり、互いに通信可能な1組のコンピュータ・デバイスまたはノードを相互接続する。ノードは、クライアントのようなコンピュータ、記憶装置およびカタログ・マネージャ、またはスイッチ、ルータ、ゲートウェイおよびその他のネットワーク・デバイスのような様々な種類の通信デバイスとすることができる。通信チャネルは、光ファイバ、同軸ケーブル、銅の撚り線対、衛星リンク、デジタル・マイクロ波無線等を含む種々の伝送媒体を用いることができる。

#### 【0062】

コンピュータ・ネットワークは、ネットワークによるノードの接続の幾何学的構成であるトポロジを有する。トポロジの種類には、二点間接続、線形バス、環状接続、星形結線、および多重接続ネットワークが含まれる。ネットワークは、これら基本的なトポロジの様々な組み合わせを用いることができる。トポロジは、物理的な設置に応じて変化する可能性がある。各ノード、即ち、クライアントまたは記憶装置を直接同じスイッチに接続する、遮断のないスイッチに基づくネ

ットワークを用いることができる。実施態様によっては、多数のクライアントおよび記憶装置を物理的ループまたはサブネットワークに接続し、これらを切り替え構造 (switching fabric) に相互接続することができる。また、本システムは、多数のスイッチを用いて接続することも可能である。

【0063】

また、ネットワークは、プロトコル、メッセージ・フォーマット、および通信ハードウェアおよびソフトウェアがネットワーク上のデバイス間で通信を行なうためのその他の規格を規定するネットワーク・アーキテクチャも有する。一般的に用いられるネットワーク・アーキテクチャは、Open Systems Interconnection reference model (オープン・システムズ相互接続基準モデル) として知られている、International Standards Organization seven-layer model (国際標準機構7レイヤ・モデル) である。7レイヤとは、アプリケーション、プレゼンテーション、セッション、トランスポート、ネットワーク、リンクおよび物理レイヤである。各機械は、これらの層の1つにおいて、同じ通信プロトコルを用いて他のあらゆる機械と通信を行なう。

【0064】

一実施形態では、リンク層は、好ましくは、パッケージがクライアントに受信される順序を保持し、無限のレイテンシの潜在的 가능성을回避するようにしたものである。したがって、適当なリンク・レイヤ・プロトコルはOC3、OC12、または更に高い帯域幅ネットワークのような、非同期転送モード(ATM)ネットワークを含む。AAL15モードで動作するATMシステムが好ましい。100Txないしギガビット(1,000Tx)の容量を有するイーサネット・ネットワークも、ソースから宛先に効率的なパケット伝送を提供する。適当なイーサネット・ネットワーク・プラットフォームは、例えば、カリフォルニア州、Santa Claraの3Comから入手可能である。ATMシステムの一例は、ペンシルベニア州、WarrendaleのFore Systems (フォア・システムズ社)、またはマサチューセッツ州、ConcordのGiga-Net (ギガネット社) から入手可能である。Fibre Channel、FD

D I、またはH I P P I ネットワークも使用可能である。異なるクライアント、カタログ・マネージャおよび記憶装置は全て、リンク・レイヤ・プロトコルを用いて通信することができる。また、このレイヤにおける通信は、各レイヤのプロトコル毎にカプセル化したデータを処理するために実行するメモリ・コピーによるオーバーヘッドを低減する。マサチューセッツ州、T y n g s b o r o の P o l y b u s S y s t e m s C o r p o r a t i o n (ポリバス・システムズ社)からの帯域幅分散ネットワーク・ファイル・システムも使用可能である。

【0065】

これまで一実施形態のコンピュータ・プラットフォームについて説明してきたが、以下では一実施形態の動作および詳細について更に追加説明を行なう。

【0066】

一実施形態では、記憶装置および当該記憶装置上に格納してあるデータを維持するためのプロセスがある。例えば、故障回復プロシージャは、ファイルの追加コピーの作成を含む場合がある。加えて、ファイルの可用性、即ち、ファイルに対するアクセスの信頼性の必要性に基づいて、ファイルを削除または追加することも可能である。最後に、保守手順によっては、記憶装置上のファイル削除、別の記憶装置へのファイル・コピー、およびシステムからの記憶装置の除去を伴う場合もある。また、ファイルは、保管したり、あるいはシステムからアーカイブ・システムに移動させることも可能である。これらのプロセスについて、図5ないし図9に関連付けて更に詳しく説明する。このようなデータ管理プロセスは、カタログ・マネージャ、別の記憶装置、またはクライアントによって行なうことができる。これらのプロセスをクライアントが行なう場合でも、カタログ・マネージャまたは記憶装置のリソースを占有することではなく、クライアントのデータ要求に対して回答するというような、他の重要性が高いタスクに用いることができる。

【0067】

図5は、記憶装置が利用不能となり、その故障を検出した後にどのようにして故障回復を行なうことができるかについて、更に詳細に記載したフローチャートである。このような故障を検出する1つの方法について、図10ないし図12に

関連付けて以下で更に詳しく説明する。要求に対して繰り返し応答し損なうことを、故障を示すために用いることもできる。このプロセスの成功は、システムないにある各セグメントのコピー数、または冗長集合内のセグメント数に依存する。コピー数をNとすると、 $N - 1$ 個の記憶装置が故障しても、システムはデータ損失なく動作する。ある記憶装置が故障した後、新たな記憶装置を代わりに設置し、失われたデータを再現するか、または失われたデータを再度作成し、残りの記憶装置全てに分配することができる。図5は、冗長データがセグメントのコピーである場合のプロセスを記載する。以下で説明する図25は、冗長情報が2つ以上のセグメントに基づく場合のプロセスを示す。

【0068】

ステップ200において、データの追加コピーを、最初にデータ、例えば、再現すべきファイルまたはソースを選択することによって行なうことができる。再現するファイルは、優先順序によって選択することができ、更に自動的または手作業のいずれかで選択することも可能である。この種の復元は、いくつかのファイルからのデータを再構成し、他のファイルからのデータが復元される前に利用可能とすることができる。ステップ202において、ソースのセグメント・テーブルを用いて、失われたデータ・セグメント、即ち、失われた記憶装置上に格納してあったデータを特定する。ステップ204において、典型的にデータを最初に捕獲するときと同様に、故障した記憶装置と交換するために新たな記憶装置が利用できない場合、失われたセグメント毎に新たな記憶装置を選択する。あるいは、交換記憶装置を選択する。ステップ206において、失われたセグメントのコピーを代替記憶装置から読み出し、選択した記憶装置に格納する。ステップ204ないし208のファイル動作は、非同期とし、各セグメント毎に別個のスレッドによって実行することができる。このような動作は、このネットワーク・アーキテクチャに備えた、多数対多数のリード／ライト機能を利用する。次に、ステップ208において、コピー動作の完了成功次に、ファイルのセグメント・テーブルを更新する。プロセスが完了した場合、カタログ・マネージャがセグメント・テーブルを維持するのであれば、ステップ209において、新たなセグメント・テーブルを用いてカタログ・マネージャを更新することができる。基のセグ

メント・テーブルが、疑似ランダム・シーケンス発生器へのシードによって表わされる場合、実際のテーブルを作成し変更しなければならない場合がある。

【0069】

このプロセスを用いた未ロード・システムに対する再入力 (repopulation) の速度および冗長性復元は、次の式で定義する。

【0070】

【数2】

$$\frac{s}{(n-1+d)(b/2)}$$

ここで、s = メガバイト単位 (MB) で示す、失われたデータのサイズ

n = 記憶装置の初期数

b = MB / 秒で表わす、記憶装置の平均帯域幅

d = MB / 秒で表わす、ユーザ・デマンドの負荷

例えば、10個の記憶装置の内1つが故障したために50ギガバイトのストレージへのアクセスが失われた場合、n = 10 記憶装置、単位帯域幅 b = 10 MB / 秒とすると、(n - 1) = 9 および (b / 2) = 5 となる。したがって、他に負荷がなければ、復元には約20分を要する。この絶対復元速度は、通常、クライアントへの変動再生負荷の逆数として減少する。例えば、50%負荷の場合、入力時間は200%増大する。呼び出されると、再分散タスクは速いレートで実行することができ、多元記憶装置チェッカボードを多数の記憶装置に切り替えるが、再入力アクティビティは適宜クライアント・ファイル・サービス要求にしたが行われる。実際の影響は、故障した記憶装置によって記憶装置の全帯域幅に多少の損失が出る程度である。復元のためのファイル選択の優先順位を決定することにより、最も重要なファイルを真っ先に復元することを保証する。

【0071】

図6は、どのようにしてデータの追加コピーを作成することができるかについて更に詳細に記載したフローチャートである。このプロセスを呼び出すと、ミッション・クリティカル (mission critical) な即ち要求の多い

データの追加データ・コピーが得られるようになる。新たなコピーに日付スタンプを与え、いつそのコピーを削除できるかを示すこともできる。選択したデータが与えられると、ステップ210において、データのセグメントを選択する。ステップ212において、各セグメントには、新たな記憶装置がランダムに割り当てられ、各記憶装置が所与のセグメントのコピーを多くても1つ有することを保証する。次に、ステップ214において、選択した記憶装置上にセグメントを格納する。このセグメントの格納完了に成功した場合、ステップ216においてデータのセグメント・テーブルを更新する。データのセグメント全てを未だコピーしていないとステップ217で判定した場合、プロセスを繰り返し、ステップ210に戻って、次のデータ・セグメントを選択する。プロセスが完了した場合、カタログ・マネージャがセグメント・テーブルを有しているのであれば、ステップ218において、新たなセグメント・テーブルによってカタログ・マネージャを更新することができる。このプロセスはセグメント全体にわたって連続的であるが、各セグメントは、別個のスレッドを用いて処理することができ、ステップ214のファイル動作は非同期とすることもできる。このような処理によって、コピーを素早く作成することを保証する。この手順により、疑似乱数発生器のシードを用いて、なおもセグメント・テーブルを表わすことができる。

#### 【0072】

図7は、データのコピーをどのようにして削除するのかについて詳細に記載したフローチャートである。このプロセスは、例えば、データがもはやさほど要求されなくなった場合に呼び出すことができる。例えば、コピー上の日付スタンプを用いて、いつデータを削除するかを示すことができる。所与のデータに対して図2に示したセグメント・テーブルを仮定すると、ステップ220において、コピー集合の1つ、即ち、テーブル内の1列を選択する。ステップ222において、列内の各セグメントを削除する。ステップ222において各セグメント毎に削除動作完了に成功した場合、ステップ224においてセグメント・テーブルを更新する。ステップ222および224は、セグメント毎に繰り返す。このプロセスは、セグメント全体を通じて連続的でもよく、または各セグメントを別個のスレッドによって処理することも可能である。プロセスが完了した場合、カタログ



・マネージャがセグメント・テーブルを維持しているのであれば、ステップ226において、カタログ・マネージャを新たなセグメント・テーブルを用いて更新することができる。

【0073】

図8は、別のアクティブな記憶装置をどのようにしてシステムから除去するのかについて記載したフローチャートである。例えば、そのファイル・システムを用いて、そのファイルのリストを識別することによって、記憶装置上で利用可能なデータを識別する。最初に、記憶装置を、新たなセグメントを書き込むために利用可能とする。このステップは、例えば、カタログ・マネージャに通知することによって、または同報通信メッセージを全クライアントに送ることによって実行することができる。各ファイルのセグメントを他の記憶装置上に再分散し、その後記憶装置をシステムから除去する。このファイル・リストが与えられると、ステップ230において、処理する次のファイルを選択する。ステップ232において、セグメント・テーブルを用いて、冗長情報を収容するセグメントを含む、記憶装置上のこのファイルのセグメント全てを識別する。ステップ234において、次に処理するセグメントを選択する。ステップ235において、残りの記憶装置からのランダム選択により、選択したセグメントに新たな記憶装置を割り当て、所与のセグメントのコピーを1つよりも多く有する記憶装置がないことを保証する。次に、ステップ236において、新たに選択した記憶装置にデータを書き込む。このライト動作の完了に成功した場合、ステップ237において、セグメント・テーブルを更新する。所与のファイルのセグメント全てを再分散したとステップ238において判定した場合、ステップ239において、セグメント・テーブルを適宜カタログ・マネージャに送ることができる。セグメントは、順次または非同期ファイル処理を用いる別個のスレッドによって処理することができる。カタログ・マネージャを更新した後に、古い記憶装置からセグメントを削除することができる。ステップ240において次のファイルがあると判定されたなら、処理は次のファイルに進む。全てのファイルを再分散し終えた場合、このプロセスは完了し、記憶装置をシステムから除去することができる。

【0074】

図9は、保管またはバックアップのためにデータをどのようにコピーすることができるかについて記載するフローチャートである。このプロセスは、利用可能な記憶装置からのデータの各セグメントの1つのコピーを、アーカイブ記憶媒体のようなバックアップ記憶システムにコピーすることを伴う。各コピー・セットおよびあらゆる冗長情報も、全ての記憶ユニットから削除することができる。このプロセスは、ステップ250において、コピー・セット、例えば、Aリストをセグメント・テーブルのコラムから選択することによって実行することができる。あるいは、前述のように他のアプリケーションによって適用される技法を用いて、各セグメントを順に読み出し、各セグメント毎に記憶装置の選択を行なうことも可能である。ステップ252において、選択したコピー・セットからの各セグメントは、その記憶装置から読み出し、記憶媒体上に格納する。各セグメントを記憶媒体にコピーするのに成功したなら、ステップ254において、全ての残りのコピー・セットまたはあらゆる冗長情報からの残りのセグメント全てを、記憶装置から削除することができる。セグメントは、順次処理することも、非同期ファイル処理を用いて別個のスレッドによって処理することも可能である。次に、ステップ256において、カタログ・マネージャを更新することができる。

#### 【0075】

図10ないし図12に関連付けて、記憶装置を監視してどのようにして可用性を判断し、故障を検出するかについてこれより説明する。記憶装置が利用可能か否かについて判定するには、いくつかの方法があり、記憶装置をポーリングすること、記憶装置からの例外を処理すること、または記憶装置が周期的にアプリケーションまたは複数のアプリケーションにその可用性について通知することを含む。一実施形態では、カタログ・マネージャ49またはその他のいずれかのクライアントに加えて、双方がシステム内においてどの記憶装置42がアクティブであるか監視し、各ファイル毎にセグメント・テーブルのカタログを維持することができる。記憶装置を監視する方法の1つを図10ないし図12に示す。システム上で利用可能な各記憶装置は、カタログ・マネージャにそれが利用可能であることを周期的に知らせるプロセスを確立する。即ち、このプロセスは、記憶装置が周期的にシステム・タイマからのタイマ割込またはイベントに応答してカウンタを

増分する、第1段階60を有する状態機械と見なすことができる。このカウンタが、100ミリ秒というような、ある所定量に到達した場合、別の状態62への遷移が発生する。状態62への遷移において、「ピング」と呼ぶ信号を記憶装置がカタログ・マネージャに送る。この信号は、1ATMセル程の小さいなメッセージとし、送信するために多くの帯域幅を用いないことが好ましい。この信号は、記憶装置の識別子、ならびに可能性として容量、効率および／または記憶装置の帯域幅可用性というような他の情報も含むことができる。次のタイマ割込またはイベント時に、カウンタをリセットし、状態60に戻る遷移が行われる。

#### 【0076】

カタログ・マネージャは、利用可能な記憶装置を追跡するとよい。この目的のために、カタログ・マネージャは記憶装置のリスト70を用いることができる。その一例を図11に示す。この記憶装置リストは、72に示す記憶装置に識別子でインデックスしたテーブルとして実施することができる。記憶装置が存在する即ち利用可能な場合、帯域幅、メモリ容量、または記憶装置のパワーに関するその他の情報を、列74において得られるようにする。列76に示すように、記憶装置からの最後の「ピング」以来のカウントもある。このカウントが、300ミリ秒というような所定量を超過した場合、記憶装置は使用不能と見なされ、続いて前述のような故障回復手順を行なうことができる。記憶装置のリスト70を維持する追跡プロセスの一例について、図12に関連付けてこれより詳細に説明する。

#### 【0077】

図12は、カタログ・マネージャがどの記憶装置が利用可能化について判定を行なう際に実行することができる追跡プロセスを記載した状態機械である。これらの状態機械の1つは、カタログ・マネージャ上のプロセスとして、各記憶装置毎に確立することも可能である。第1状態80は、待ち状態であり、記憶装置のリスト70内の記憶装置に対するカウント値76を、周期的なタイマ割込に応答して、当該記憶装置に対して増分する。記憶装置から「ピング」を受信した場合、状態82への遷移が起こる。状態82において、リスト70内のこの記憶装置の存在を検証する。記憶装置がリスト70内にある場合、この記憶装置のカウン

ト76をリセットし、この記憶装置に関する情報を更新することができ、そして状態80に戻る遷移が行われる。記憶装置がリストにない場合、これをリストに追加し、カウントをリセットし、状態80に戻る遷移が行われる。所与の増分の後、記憶装置のカウントが、300ミリ秒というような所定のタイム・アウト値よりも大きくなった場合、故障回復手順を実行する。即ち、リスト70から記憶装置を除去し、状態84において、フォールト・トレラント手順を実行する。記憶装置からの「ピング」をカタログ・マネージャが受信し、その記憶装置が対応する追跡プロセスを有していない場合、カタログ・マネージャは記憶装置をリストに追加し、この記憶装置に対して追跡プロセスを作成する。

【0078】

カタログ・マネージャ49を有することに加えて、システムは、アセット・マネージャと呼ぶデータベースも含むことができる。アセット・マネージャは、各ファイル毎のインデックスのように、システムにおいて利用可能なメディア・ソースに関する種々の情報を格納する。カタログ・マネージャおよびアセット・マネージャは組み合わせることも可能である。アセット・マネージャに格納する有用な一種の情報は、図13に示すテーブルであり、これは、ソース識別子および、米国特許第5,267,351号に示されるような、当該ソース内における範囲に基づいて、等価データ・ファイルに関係付けるものである。ソース識別子は、データの原ソースの指示であり、アナログ・ソースとすることができる。これに対して実際に利用可能なデータは、記憶装置上に格納してある当該ソースのデジタル化したコピーである。即ち、テーブルは、ソース識別子のエントリ、ソース識別子102内の範囲、および、データ・ファイルのリストのような、当該ソースからの等価なデータの指示104を有する。リスト104は、ソースに対するデータ・ファイルの1つを特定し次いで当該ファイルに対するセグメント・テーブルにアクセスし、データのセグメントが種々の記憶装置上のどこに分散されたかについて判定するために用いることができる。図2Aのセグメント・テーブル90Aは、106および108に示すように、図13のこのリストに組み込むこともできる。図2Bのセグメント・テーブル90Bも同様に、リスト104に組み込むことができる。このようなデータ間の等価性は、あらゆるアプリケー

ション・プログラムによっても維持することができる。

【0079】

カタログ・マネージャは、種々の記憶装置上にどのようにデータを分散したかを監視するデータベースであるので、フォールト・トレランスおよび可用性を高め、ボトルネックのなるその確度を低下するように設計して当然である。したがって、カタログ・マネージャは、従来の分散データベース管理技法を用いて実現する。また、Mrathon Technologies（マラソン・テクノロジー社）、Tandem Computers（タンデム・コンピュータズ社）、Stratus（ストラタス）、およびTexas Micro, Inc.（テキサス・マイクロ社）からのもののような、非常に可用性が高い機械を用いて、カタログ・マネージャを実現することも可能である。また、別個のクライアント・アプリケーションが用いるカタログ・マネージャもいくつかあり得る。あるいは、各クライアント・アプリケーションは、標準的な技法を用いて多数のデータ・コピー間の一貫性を維持し、それ自体のカタログのコピーをローカルに維持することも可能である。このように、カタログ・マネージャは、故障の中心点ではない。また、クライアントは、それ自体のカタログ・マネージャとして作用することも可能である。また、カタログは、データとして扱うことも可能であり、そのセグメントおよび冗長データを記憶装置間でランダムに分散する。各クライアントは、各カタログ毎に、セグメント・テーブル、または、セグメント・テーブルを表わす乱数発生器のシードを有することもできる。

【0080】

以上、どのようにしてデータを捕獲し記憶装置に格納することができるか、そしてどのようにして記憶装置上のデータ格納を管理することができるかについて説明したので、これよりオーサリングおよび再生を行なうクライアント・アプリケーションについて、図4および図15に関連付けて更に詳しく説明する。

【0081】

マルチメディア・データのオーサリング、処理および表示を行なうには、いくつかの種類のシステムを用いることができる。これらのシステムは、データを修正し、異なるデータの組み合わせを定義し、新たなデータを作成し、データをユ

ーザに表示する際に用いることができる。これらの種類のシステムを実現する種々の技法が、当技術分野では公知である。

【0082】

マルチメディア・オーサリング、処理および再生システムは、典型的に、マルチメディア組成 (multimedia composition) を表わすデータ構造を有する。このデータ構造は、最終的に、一意の識別子またはファイル名のようなソース素材の識別子、そして可能性としてクリップを定義するソース素材内の時間的範囲を用いて、デジタル化ビデオまたはオーディオのような、ソース素材のクリップを引用する。識別子は、等価データ・ファイルのリストと共に用いて、ソース素材のファイル名を特定するものであれば、いずれの形式でもよい。インデックスを用いて、ソース内の時間範囲を、対応するファイル内のバイト範囲に変換することも可能である。このバイト範囲を、ファイルのセグメント・テーブルと共に用い、必要とされるセグメントおよびデータを検索する記憶装置を特定することができる。

【0083】

図14は、マルチメディア組成の一部を表わすために用いることができる、リスト構造の一例を示す。図14に示す例では、数個のクリップ260があり、その各々が、262で示す、ソース識別子への引用、および264で示すように、ソース内の範囲を含む。一般に、時間組成において、メディアの各トラック毎に、このようなリストがあるとよい。組成を表わすために使用可能な種々のデータ構造がある。リスト構造に加えて、1993年10月28日に公開されたPCT公開出願WO93/21636には、更に複雑な構造が示されている。他のマルチメディア組成の代表例には、Avid Technology, Inc. からのOpen Media Framework Interchange Specification、Multimedia Task ForceからのAdvanced Authoring Format (AAF)、MicrosoftからのDirect Show、およびPCT公開WO96/26600に示されているような、Apple ComputerからのBentoが含まれる。

## 【0084】

先に述べ、マルチメディア・プログラムを表わすために用いられるデータ構造は、多数の形式のデータを用い、これらを同期させて表示することができる。最も一般的な例は、モーション・ビデオ（多くの場合、2つ以上のストリーム即ちトラック）および付随するオーディオ（多くの場合、4つ以上のストリーム即ちトラック）を含む、テレビジョン番組または映画の製作である。図15に示すように、クライアント・コンピュータは、主メモリ内に割り当てたメモリ・バッファ294の対応するセット290を有することができる。各バッファは、「直列化」バッファとして実現することができる。言い換えると、クライアントは、記憶装置から受信したデータを、これら独立してアクセス可能な部分に挿入し、バッファ集合から順次読み出す。要求を数個の記憶装置に送ることができ、データは同じストリームに対して異なる時点で受信することができるので、バッファは、書き込み時には、連続的に満たされないこともあるが、読み出しは順次行われ表示される。図15では、バッファ内の充填は、バッファ内におけるデータの存在によって示す。293および295で示すように、空のバッファはいずれも、いずれの時点でも充填することができる。しかしながら、各バッファ集合は、現リード・ロケーション291を有し、ここからデータを読み出し、かつ297に示すように、時間の進展に連れて前進する。これらのバッファの部分集合292, 296を各データ・ストリームに割り当てることも可能である。

## 【0085】

バッファ集合内の各バッファは、固定数のデータ・セグメントに対応するサイズを有し、セグメント・サイズとは、記憶装置上に格納したセグメントのサイズである。オーディオ・データ・ストリーム292当たり数個、例えば、4つのオーディオ・バッファがあるとした場合、各バッファは、数個、例えば、4つのセグメントを収容することができる。同様に、各ビデオ・ストリーム296は、数個、例えば、4つのバッファを有することができ、その各々は、数個、例えば、4つのセグメントを収容する。バッファの各々は、独立してアクセス可能な部分298に分割することができる。部分298のサイズは、ネットワークを通じた転送を予定しているデータ・パケットのサイズに対応する。

## 【0086】

ビデオおよびオーディオ・データは異なるデータ・ファイルに格納することができ、しかも任意に組み合わせることができるので、クライアント側でこれら異なるストリームのデータ要求を効率的に管理すれば、一層の性能向上が得られる。例えば、クライアント・アプリケーションは、データを読み取ることができるストリームを識別し、次いで、データがあれば読み取るべきデータ量を決定することができる。この種のリード動作管理を行なうプロセスが、米国特許第5,045,940号に示されている。一般に、クライアントは、表示に利用可能なデータ量が最少のストリームはどれであるかについて判定を行なう。ある量のデータを効率的に読み出すために当該ストリームのために十分な量のバッファ空間がバッファ集合にある場合、そのデータを要求する。一般に、選択したストリームに対してメモリ内で利用可能な空間が十分に広く1ネットワーク伝送データ単位を保持できる場合、データは効率的に読み取られる。あるストリームのデータを要求すると判定した場合、各セグメントが格納されている記憶装置から選択した記憶装置から、データの各セグメントを要求する。

## 【0087】

ある組成をデータ要求に変換し、当該データを表示するプロセスについて、図16に関連付けてこれより説明する。記憶装置からどのファイルを要求するか知るために、クライアント・システム上で実行するアプリケーション・プログラムは、図14に示すような組成を表わすデータ構造を、図16のステップにおいてこれらのファイル内におけるファイル名および範囲に変換することができる。例えば、各ソース識別子およびそのソース内の範囲毎に、要求をアセット・マネージャに送ることができる。これに応答して、アセット・マネージャは、受信したソース識別子および範囲に対応する等価なメディアを収容するファイルのファイル名を戻すことができる。ファイルのセグメント・テーブルおよび利用可能な記憶装置のリストも、カタログ・マネージャとなることができる。

## 【0088】

クライアントが特定のデータ・ストリームに対してデータ・セグメントを要求する場合、クライアントはステップ272において、要求したセグメントのため



に記憶装置を選択する。各セグメントをコピーすることによって冗長性を備えた一実施形態におけるこの選択について、図17および図18に関連付けて以下で更に詳しく説明する。一般に、キュー48（図1）が最も短い記憶装置を選択すればよい。次にクライアントは、ステップ274ないし278において、選択した記憶装置から、そのセグメントに対するデータを読み出す。ステップ274は、プレリード・ステップとして理解するとよく、ここで、クライアントは記憶装置に要求を送り、所望のデータを、不揮発性ストレージからより高速な典型的に揮発性のストレージに読み出す。記憶装置に対する要求は、要求を行ってからクライアントにおいて要求データを受信しなければならないまで、即ち、期限までに要する時間がどれ位かかるかの指示を含むことができる。プレリード要求が受け入れられると、クライアントはステップ276において待機する。要求は記憶装置のキュー48に置かれ、次いで期限を用いて以下に説明するように要求の優先順位を決める。ステップ278において、記憶装置のバッファにおいてデータが利用可能となった後、記憶装置からデータを転送する。このステップは、ネットワークの使用をスケジューリングしネットワーク利用効率を最大に高めることを含むとよい。受信したデータは、クライアントにおける適切なバッファに格納し、最終的にステップ280において処理し表示する。記憶装置においてセグメントが失われていた場合、冗長情報を用いてセグメントを再現する。

#### 【0089】

プレリード要求を開始する方法はいくつかあり、ステップ274における記憶装置の選択、およびステップ278におけるデータ転送を含む。例えば、マサチューセッツ州、TewksburyのAvid Technology, Inc.からのMedia Composerは、ユーザがクリップ数または時間量のいずれかをルック・アヘッド値（look-ahead value）として設定させ、組成においてどの位先でアプリケーションがデータ・リード要求を開始すべきかを示す。テレビジョン放送説明のためのプログラム・スケジュールもこの目的のために使用可能である。このような情報は、記憶装置の選択およびプレリード要求を開始するために用いることができる。このようなプレリードは、1995年9月9日公開のヨーロッパ特許出願第0674414A2号に示されるよ

うに、バッファ290（図15）内にバッファ空間が得られない場合でも実行可能である。バッファ290（図15）における利用可能な空間量を用いて、ステップ278（図16）においてデータ転送を開始するか、あるいはプレリード（ステップ274）およびデータ転送（ステップ278）を開始することができる。

#### 【0090】

利用可能な記憶装置全てを洗いざらい探索する必要なく、クライアントが、最も短い要求キューを有する記憶装置の適切な推定を行なうことを可能にする1つのプロセスについて、図17および図18に関連付けてこれより説明する。最初に、クライアントは、ステップ330において、スレシホールドE1を有する要求を第1記憶装置に送る。スレシホールドE1は、要求に答えなければならない時間の推定値を示す。この推定値は、時間値、4というような記憶装置のディスク・キュー内の要求数、またはその他の尺度として表わすことができる。このスレシホールドの意味は、例えば、記憶装置が特定の時間制限の内に要求に答えることができる場合に、記憶装置によって要求が受け入れられなければならないということである。クライアントは、ステップ332において、記憶装置から回答を受信する。回答は、要求が受け入れられ記憶装置のディスク・キューに入力されたか、または要求が拒絶されたかを示す。これは、ステップ334において判定される。要求が受け入れられた場合、ステップ336において、記憶装置のバッファにおいてデータが利用可能となる時点の推定値がクライアントには与えられる。例えば、要求したセグメントのデータが既にバッファ内にある場合、記憶装置はデータが直ちに利用可能であることを示す。次いで、クライアントは、推定時間が経過して少し後に、データの転送を要求する時刻になるまで待てばよい。要求が拒絶された場合、ディスク・キューのエントリ数の実際のサイズのような、記憶装置が実際に要する可能性がある時間量の推定値が、記憶装置から戻される。ステップ340において、この実際の推定値を値Kに加算し、スレシホールドE2を得る。値Kは、ディスク・キューのエントリ数を表わす場合は、2となる。スレシホールドE1および値Kはユーザが定義可能である。ステップ342において、スレシホールドE2を示す要求を第2記憶装置に送る。次に、ステップ344に

において、クライアントは、ステップ332において受信した回答と同様に、回答を受信する。要求が受け入れられたことを回答が示すと346において判定した場合、クライアントは、ステップ336に示すように、第2記憶装置においてデータが利用可能となる時点の推定値を有し、その後クライアントは、データ転送をスケジュールするために待つことができる。それ以外の場合、ステップ348において、無条件要求、即ち、大きなスレシホールドを有する要求を第1記憶装置に送る。次に、ステップ350において承認を受信し、ステップ336において示したように、記憶装置のバッファにおいてデータが利用可能となる時点の推定値を示す。

#### 【0091】

一方、記憶装置は、要求を受信したときには、それがクライアントによって選択された第1記憶装置なのかまたは第2記憶装置なのかはわからない。逆に、記憶装置は、ステップ360に示すように、単に要求を受信するだけである。ステップ362において、要求内で示されるスレシホールドは、記憶装置自体の、クライアントが待たなければならない時間の推定値と比較される。例えば、記憶装置のディスク・キューのサイズを指定されたスレシホールドと比較する。供給内のスレシホールドが、記憶装置の推定値よりも大きい場合、ディスク・キューに要求を入力し、ステップ364において、記憶装置のバッファにおいてデータが利用可能となる時点の推定値を決定する。この推定値は、例えば、ディスク・アクセス速度、ディスク・キューの長さ、そして可能性として最新のパフォーマンスの移動平均に基づいて決定することができる。ステップ366において、承認をクライアントに送る。これは、記憶装置のバッファにおけるデータ可能性の推定時点を含む。それ以外の場合、ステップ368において拒絶を送り、ディスク・キューの実際の長さのような、この推定値を示す。

#### 【0092】

記憶装置は、記憶装置上のどのバッファに、どのセグメントがあるのか追跡することも可能である。セグメント・データは、記憶媒体からいずれかの空きバッファまたは最近最も使用されていないセグメントが占有するバッファに読み出すことができる。このように、あるセグメントに対するデータは、当該セグメント

を2回要求すれば、あるバッファにおいて直ちに利用可能とすることができる。

【0093】

代替案として、クライアントは別の方法を用いて、データを検索する記憶装置を選択することも可能である。これについては、以下で論ずる。要求を送った後、クライアントは、記憶装置から、その要求が記憶装置のディスク・キュー内に入ったことを示す承認を受信する。記憶装置のバッファにおいてデータが利用可能となる時点の推定値を受信する代わりに、クライアントは、記憶装置が要求データを記憶装置の指定バッファ・メモリに読み込んだことを示すレディ信号を受信するまで待つこともできる。この待ち時間中、クライアントは、別のデータ・セグメントの要求を発行したり、データを表示したり、データを処理するというような、別のタスクを実行することができる。この代替案に伴う1つの問題は、クライアントが未要請メッセージ、即ち、記憶装置からのレディ信号を受け入れ、それに応答して、クライアントがコンテキストを変更し、メッセージを処理することである。クライアントは、他の動作を実行していて使用中の可能性もある。このプロセスは記憶装置のバッファにおいてデータが利用可能となる時点の更に精度が高い推定値を与えるが、コンテキストを変更し着信メッセージを処理する可能性のために、直ちにクライアントにおける複雑化を招くことになる。

【0094】

セグメント・テーブルが各セグメントのコピーを追跡する場合に、あるファイルのためにセグメント・テーブルから記憶装置を選択する方法は、他にもいくつかある。例えば、クライアントがファイル・リード要求を行なっている場合、クライアントは、問題のファイルのために、「A」リストまたは「B」リストのいずれかからランダムに抽出することができる。あるいは、クライアントは、その現在未処理の要求全て、即ち、送られたが未だ実行していない要求を見直し、AおよびBリスト上の記憶装置の内、最も少ない未処理の要求を有する記憶装置を、ファイルに抽出することができる。この選択方法は、クライアントがそれ自体の未処理要求と競合する可能性を低下させることができ、要求を全ての記憶装置全体に一層等しく拡散するようになる。あるいは、未処理の要求を調べる代わりに、クライアントは、その最新要求の履歴、例えば最後の「n」の要求を調べ、

次の要求のために、これまで最も使用されていないセグメントに対するAリストおよびBリストからいずれかの記憶装置を抽出することができる。この選択方法は、要求を全ての記憶装置全体に一層等しく拡散するようになり、更に特定の記憶装置に要求が集中するのを回避することができる。また、クライアントは、各記憶装置から、そのディスク・キューの長さの尺度を要求することもできる。クライアントは、最もディスク・キューが短い記憶装置に要求を発行することができる。別の可能性として、クライアントは、2つの記憶装置に要求を送り、最終的に一方のみからデータを受信することもできる。この方法をローカル・エリア・ネットワーク上で用いると、クライアントは未使用の要求を取り消すことができる。ワイド・エリア・ネットワーク上では、最終的に選択された記憶装置は、他の記憶装置において未使用の要求を取り消すことができる。

【0095】

記憶装置は、多数のアプリケーションから多数の要求を受信する可能性もある。多数のアプリケーションからの要求を管理し、最も重要な要求を最初に処理することを保証するために、各記憶装置毎にキュー48（図1）を維持する。キューは、システムの複雑度に応じて、いくつかの部分に維持することも可能である。即ち、記憶装置は、ディスク・アクセスのためおよびネットワーク転送のために異なるキューを維持することができる。キューは、例えば、再生から放送のために特定の時間期限があるデータを用いる時間厳守アプリケーションからの要求を、キャプチャ・システム、オーサリング・ツールまたはサービスおよび保守アプリケーションのようなその他のアプリケーションからの要求から、分離することも可能である。更に、格納要求は、オーサリング・ツールからの要求、ならびにサービスおよび保守プログラムからの要求から分けることができる。オーサリング・ツールからの要求は、更にサービスおよび保守要求から分けることができる。

【0096】

図19は、ディスク・キュー300およびネットワーク・キュー320を利用するキュー48の一実施形態を示す。ディスク・キューは4つのサブキュー302, 304, 306および308を有し、それぞれ、再生、キャプチャ、オーA

リング、ならびにサービスおよび保守クライアント・プログラムに1つずつである。同様に、ネットワーク・キュー320も4つのサブキュー322, 324, 326および328を有する。各キューは、1つ以上のエントリ310を含み、その各々は、要求を行なったクライアントおよび要求された動作を示す要求フィールド312、要求の優先度を示す優先度フィールド314、ならびに要求に関連するバッファを示すバッファ・フィールド316を含む。要求の優先度の指示は、デッドライン、タイム・スタンプ、クライアントにおいて利用可能なメモリ量の指示、またはクライアントにおいて現在利用可能なデータ量の指示とすることができる。記憶装置に優先度スケジューリング機構があれば、使用する優先度スタンプの種類を指示することができよう。

#### 【0097】

優先度値は多くの方法で生成することができる。オーサリングまたは再生システムに対する優先度値は、一般に、アプリケーションが要求データを受信しなければならない時点の尺度である。例えば、リード動作では、アプリケーションは、データがなくなる目に、再生に利用できるデータがどの位あるかについて（ミリ秒単位またはフレーム単位またはバイト単位で）報告することができる。キャプチャ・システムに対する優先度指示は、一般に、クライアントがそのバッファから記憶装置にデータを転送しなければならない時点の尺度である。例えば、ライト動作では、アプリケーションは、バッファが溢れる前に満たすことができる空のバッファ空間がどれ位あるかについて（ミリ秒単位、フレーム単位またはバイト単位で）報告することができる。ミリ秒を尺度単位として用いた場合、システムは、絶対的な時間クロックを有し、これをキュー49における要求の順番を決定するための基礎として用いることができ、全てのアプリケーションおよび記憶装置をこの絶対時間クロックに同期させることができる。この同期が実用的でない場合、アプリケーションは、当該アプリケーションに関係し、要求を行って渡されてから、要求データがクライアントによって受信されるまでどれ位の時間を要するかを示す時間を用いることも可能である。通信レイテンシが低いと仮定すると、記憶装置は、この相対時間を、記憶装置と一貫性のある絶対時間に変換することも可能である。

## 【0098】

記憶装置は、そのサブキュー301ないし308において、要求をその優先度の順に処理する。即ち、最初に最も優先度が高いキューにある要求をそれらの優先度値の順に処理し、次いで、次第に優先度が低くなるキューにおける要求を処理する。各要求毎に、記憶装置は、ディスクと要求によって示されるバッファとの間でデータを転送する。リード要求では、要求を処理した後、ディスク・キューからネットワーク・キューに要求を移管する。ライト要求では、ライト動作の完了に成功した後、ディスク・キューから要求を除去する。

## 【0099】

以下で更に詳しく説明する一実施形態では、記憶装置はネットワーク・キューを用いて、これらの転送をスケジュールするプロセスにおいて、ネットワーク転送の優先順位を決定する。この実施形態では、クライアントは、ネットワークを通じてデータの転送を要求する。記憶装置がこのような要求を2つほぼ同時に受信した場合、記憶装置は、そのネットワーク・キューにおいて優先度が高い方の要求を処理する。リード要求では、要求を処理した後、ネットワーク・キューから要求を除去する。ライト要求では、転送完了が成功した後、要求をネットワーク・キューからディスク・キューに移管する。優先度は、空きバッファの可用性に依存する。ネットワーク・キュー内にある要求を処理する時間が過ぎた場合、要求を欠落させ、クライアントはもはや動作していないか、またはネットワーク転送を時間内に要求しなかったことを示す。

## 【0100】

記憶装置とクライアントとのコンピュータ・ネットワークを通じたデータ転送をスケジュールし、効率を向上させることも可能である。即ち、データ転送のスケジュールリングにより、コンピュータ・ネットワークの帯域幅利用度を高める。このようなネットワーク使用のスケジュールリングは、特にクライアントとスイッチとの間のリンクの帯域幅が、記憶装置とスイッチとの間のリンクの帯域幅と同程度の大きさである場合に実行するとよい。即ち、記憶装置がデータを送り、クライアントがそれらそれぞれのネットワーク接続のリンク速度でデータを受信する場合、データはネットワーク・スイッチにおいて蓄積したり、他の大幅な遅れ

が生ずる可能性は低くなる。

【0101】

このようなネットワークの利用度を高めるために、いずれの所与の時点においても、各クライアントに1つの記憶装置のみからデータを受信させ、各記憶装置に1つのクライアントのみにデータを送らせる機構を備えるとよい。例えば、各クライアントがトークンを1つのみ有するようにするとよい。クライアントは、このトークンを1つの記憶装置のみに送り、選択したセグメントに対するデータの転送を要求する。トークンは、クライアントがデータを受信しなければならないデッドライン、即ち、優先度の尺度および指定セグメントを示すことができる。各記憶装置は一度に、それがトークンを受信した1つのクライアントだけにデータを送る。記憶装置は、一度に1つのトークンのみを受信する。データを転送した後、記憶装置はトークンも返す。

【0102】

図20および図21に関連付けて他のネットワーク・スケジューリング・プロセスについてこれより説明する。このプロセスは、同様の結果を与えるが、トークンを使用しない。むしろ、クライアントが記憶装置との通信チャネルを要求し、セグメント、およびクライアントが転送が行われるのを待つ時間量E3を指定する。また、クライアントは、クライアントがデータを受信しなければならない新たな時間期限をセグメントに対して指定することもできる。

【0103】

これより図20を参照して、ネットワークを通じてデータを転送するクライアント・プロセスについて説明する。組成の再生中のいずれかの時点において、各バッファは、それに関連するデータ・セグメント、および連続再生のためにバッファ内においてデータが利用可能でなければならない時点を有する。当技術分野では公知であるが、アプリケーションは、再生プロセスの間、バッファの各々をセグメントに関連付ける。図17および図18に関連付けて先に示したように、クライアントが準備した各セグメントには、記憶装置においてデータが利用可能となる推定時間が関連付けられている。したがって、クライアントは、バッファをそれらの時間期限によって、および要求データが記憶装置のバッファ内におい



て利用可能であることが求められているか否かによって順序付けることができる。この順序付けは、クライアントが、ステップ500においてデータを転送するために次のバッファを選択するために用いることができる。クライアントは、ステップ502において、記憶装置との通信チャンネルを要求し、待ち時間E3を指定する。クライアントがデータを早急に必要としない場合、またはクライアントが他の動作をより効率的に実行できる場合、この値E3は、例えば、100ミリ秒と短くてよい。この値E3は、クライアントが早急にデータを必要とする場合には、例えば、そのバッファの1つでデータが無くならないように長くすることも可能である。ステップ504において、クライアントは記憶装置から回答を受信する。記憶装置が要求を拒絶したとステップ506において判定した場合、ステップ508において改訂推定時間をメッセージと共に受信する。ステップ510において、この改訂推定時間を用いて、バッファを選択するバッファ・リストを更新することができる。処理はステップ500に戻り、別のバッファを選択する。セグメントが、以前に選択したセグメントと同じ記憶装置上にあるバッファは、恐らく選択すべきでない。記憶装置が別の方法で要求を受け入れた場合、最終的にステップ518でデータを受信する。

#### 【0104】

記憶装置の観点からの処理について、図21に関連付けてこれより説明する。ステップ520において、記憶装置は、待ち時間E3を示す要求をクライアントから受信する。この記憶装置のバッファにはデータが未だ利用可能でないとステップ522において判定した場合、ステップ524において記憶装置は要求を拒絶し、改訂推定時間を計算し、これをクライアントに送る。逆に、データが利用可能であり、記憶装置のネットワーク接続が使用中でないとステップ526において判定した場合、ステップ528においてクライアントは「アクティブ・クライアント」となり、記憶装置によって通信チャンネルが付与され、データの転送が可能となる。記憶装置のネットワーク接続が他のクライアントにデータを転送するために稼働している場合、記憶装置は「待機クライアント」からの要求を維持し、「アクティブ・クライアント」に対するデータ転送が完了した後に、「待機クライアント」にデータを転送する。現クライアントが「待機クライアント」

か否かについて判定を行なうために、ステップ530において、記憶装置は、ネットワーク・キューにおいてデッドラインが速い要求の数に、各要求毎のネットワーク伝送時間を乗算した値に基づいて、転送を行なうことができるまでの時間を推定する。計算した利用可能推定時間が待ち時間E3よりも大きいとステップ532で判定した場合、クライアントはそんなに長く待ちたくないことを示し、ステップ524において要求を拒絶する。また、この要求の指定優先度が、いずれの現待機クライアントの優先度よりも低いとステップ534において判定した場合、ステップ524において要求を拒絶する。それ以外の場合、ステップ536において、いずれの現待機クライアントからの要求も拒絶し、この新たなクライアントを現待機クライアントとして指定する。アクティブ・クライアントへの転送が完了すると、待機クライアントがアクティブ・クライアントとなり、データを転送する。

#### 【0105】

クライアントから記憶装置にデータを転送するためには、同様のプロセスを用いて、ネットワーク転送をスケジュールし、記憶装置内のバッファから不揮発性ストレージにデータを転送することができる。クライアントの観点から、このプロセスについて図22に関連付けてこれより説明する。このプロセスは、図3におけるステップ124および126を実行するために用いることができる。

#### 【0106】

クライアントがデータをそのバッファ集合内の任意の点に入力することができるリード・プロセスとは異なり、記憶装置に転送するデータは、典型的に、キャプチャ・システムが使用するバッファ集合からのリード・ポインタから来る。

#### 【0107】

キャプチャ・システムは、典型的に、1つ以上のビデオ情報ストリーム、および1つ以上のオーディオ情報ストリームを生成する。したがって、キャプチャ・システムは、ストリーム内の空きバッファ空間量に応じて、データ・ストリームの1つを選択し、捕獲したデータを受信することができる。この選択したストリームの現リード・ポインタにおけるバッファを、ステップ600において選択する。次に、ステップ602においてライト要求を記憶装置に送る。要求は、セグ

メントの識別子、時間期限またはその他の優先度値、およびクライアントが待とうとする時間量を示すスレシホールドE4を含む。時間期限は、記憶装置がネットワーク転送要求の優先度を決定するために用いる。スレシホールドE4は、先に論じたスレシホールドE3と同様、クライアントが用い、それ自体の動作を効率的にスケジュールすることを可能にする。クライアントは、要求を記憶装置に送った後、最終的にステップ604において回答を受信する。ライト要求が拒絶されたことを要求が示すとステップ606において判断した場合、回答は、ステップ607においてデータを受信するために記憶装置が利用可能となるまでの推定時間を含む。この推定時間は、クライアントが他の動作をスケジュールするために用いることができる。記憶装置がデータを書き込む要求を受け入れた場合、ステップ608において、クライアントは、データのセグメントの一部を記憶装置に送る。ステップ610において、ライト要求が成功したか否かを示す回答をステップ610で受信することができる。これは、ステップ612において分析する。ステップ614において、不良の場合には、復元プロセスを伴う場合もある。それ以外の場合、ステップ616に示すように、プロセスは完了する。

#### 【0108】

記憶装置の観点から、記憶装置はステップ620においてクライアントからのライト要求を受信する。要求は、時間期限またはその他の優先度スタンプを含み、これを用いて要求をネットワーク・キューに入力する。次に、ステップ622において、記憶装置はデータを受信するためにバッファが利用可能か否かについて判定を行なう。利用可能なバッファがないというありそうもない場合には、ステップ624において要求が拒絶される場合もある。それ以外の場合に、ステップ626において要求をネットワーク・キューに入力し、データ、その優先度スタンプ、および転送に関するその他の情報を受信するためにバッファを割り当てたことを示す。次に、記憶装置は、ステップ628において、ネットワーク接続が使用中か否かについて判定を行なう。ネットワーク接続が使用中でない場合、記憶装置はステップ630において要求を受け入れ、その旨のメッセージをクライアントに送る。すると、クライアントはデータを転送し、これがステップ632において記憶装置によって受信され、指定のバッファに入力される。指定のバ

ッファが現在満杯であるとステップ634において判定した場合、ステップ636において、適切な優先度スタンプと共にバッファをディスク・キューに入力する。記憶装置のディスク・キューの処理によって、最終的にデータはバッファから永続ストレージに転送される。それ以外の場合、ステップ638に示すように、クライアントがバッファを満たすのに十分なデータを送るまで、記憶装置は待機する。

【0109】

記憶装置のネットワーク接続が使用中であるとステップ628において判定した場合、記憶装置はステップ640において、記憶装置のネットワーク接続が利用可能となるまでの推定時間を計算する。この計算時間が、ステップ642に判定した指定の待ち時間E4よりも大きい場合、ステップ643において、記憶装置の利用可能時間の推定値によって、要求を拒絶する。記憶装置が、クライアントによって示された待ち時間E4以内にデータを転送することができる場合、ステップ644において、記憶装置は、要求の優先度を、現在待機中のあらゆるクライアントの要求の優先度と比較する。この要求の優先度が、現在待機中のクライアントの要求よりも低い場合、この要求を拒絶する。それ以外の場合、現在待機中のクライアントからの要求を拒絶し、この新たな要求を次の要求とし、ステップ646において処理する。

【0110】

冗長情報を2つ以上のセグメントから作成する場合に用いる追加実施形態について、これより図24および図25に関連付けて説明する。

【0111】

まず図24を参照して、数個の記憶装置にわたってランダムに分散して冗長情報と共にデータ・セグメントを格納するプロセス例について詳細に説明する。このプロセスは、図3に関連付けて先に説明したプロセスと全体的に類似している。最初に、ステップ700において、キャプチャ・システムはセグメント・テーブル90B（図2B）を作成する。典型的に、各画像を、捕獲するデータ・ストリームにおけるオフセットにマップする画像インデックスも作成する。インデックス化した画像は、例えば、ビデオのフィールドまたはフレームに対応すること

ができる。インデックスは、オーディオのような他の種類のデータに対する、時間期間のような、別のサンプル境界を基準とすることができる。また、キャプチャ・システムは、前述のような、利用可能な記憶装置のリストも得る。また、キャプチャ・システムは、冗長セット・サイズの指示を、利用可能な記憶装置のリストに基づいて自動的に、またはユーザから受け取る。一般に、冗長セット・サイズは、利用可能な記憶装置の数未満でなければならない、非常に小さい部分集合であってもよい。また、カウンタを用いて、どのセグメントが所与の冗長セット内にあるか追跡する。このカウンタは、ステップ700においてゼロにリセットする。排他的ORメモリも使い、全て二進無アサート値、例えば、「0」にリセットする。

#### 【0112】

次に、ステップ720で、データ・セグメントをキャプチャ・システムにおいて作成する。このセグメントに適切なサイズについては、図3の説明との関係で先に論じた。また、ステップ720ではカウンタの増分も行なう。

#### 【0113】

ステップ722において、現セグメントは、排他的ORメモリに既に格納したあらゆるセグメントの排他的ORとしてローカルに格納する。ステップ724において、記憶装置をこのセグメントのために選択する。セグメントのための記憶装置の選択は、ランダムまたは疑似ランダムである。この選択は、以前のいずれの冗長セットに対して行なった選択にも独立とすることができる。しかしながら、選択は、冗長セット内の各セグメントが異なる記憶装置上に格納されていることを保証しなければならない。図3の説明との関係で先に論じたように、各ファイルは、利用可能な記憶装置の部分集合のみを使用することができる。

#### 【0114】

セグメントに記憶装置を選択した後、ステップ726においてこのセグメントを記憶装置に送り格納する。次に、ステップ728において、キャプチャ・システムは、記憶装置がセグメントの格納完了を承認するまで待つ。データを捕獲しながらリアル・タイムで格納しなければならない場合、ステップ726におけるデータ転送は、先に論じたようにリード動作と同様、2段階で行なえばよい。デ

ータを記憶装置に格納するのに成功した後、ステップ730においてキャプチャ・システムはセグメント・テーブル90Bを更新する。

【0115】

カウンタが現在冗長セットのサイズに等しいとステップ732において判定した場合、ローカル排他的ORメモリの内容は冗長情報である。この場合、この冗長情報を記憶装置上に格納する。即ち、ステップ734においてカウンタをリセットする。ステップ736において、冗長情報のために記憶装置を選択する。ステップ738において、冗長情報を選択した記憶装置に送る。次いで、ステップ740において、キャプチャ・システムは格納成功の承認を待つ。ステップ742において、セグメント・テーブルを更新することができる。

【0116】

捕獲が完了したとステップ744において判断した場合、プロセスは終了する。この時点で、ステップ745において、排他的ORメモリに格納されているあらゆる冗長情報は、ステップ734ないし742と同様の手順を用いて、記憶装置に格納しなければならない。次に、ステップ746において、更新したセグメント・テーブルをカタログ・マネージャに送る。ステップ732においてカウンタが冗長セットのサイズに等しくない場合、また捕獲が完了していないとステップ744において判断した場合、プロセスは継続し、ステップ720においてデータの次のセグメントを作成し、カウンタを増分する。

【0117】

図5に関連付けて先に論じたように、冗長情報は、記憶装置の1つが故障した場合に、データを復元することを可能にする。図25は、冗長情報が、2つ以上のセグメントを収容する冗長セットに基づく場合に、このような故障回復を行なうプロセスを示す。図5に示すように、復元するファイルをステップ750において選択する。ステップ752において、そのファイルの失われたセグメントを全て特定する。次に、ステップ754において、失われたセグメントを収容する冗長セットを読み取る。このステップは、当該セットにおけるセグメントの排他的ORによって作成したセットに対する冗長情報を読み出し、次いで冗長セットの残りのセグメントを読み出すことからなる。次に、ステップ756において残

りのセグメントおよび冗長情報の排他的ORを計算し、失われたセグメントを再構成する。図5のステップ204と同様に、次に、ステップ758において、失われたセグメントを再構成したものの毎に記憶装置を選択する。選択した記憶装置に、失われたセグメントを再構成したものを格納する。ステップ760において、格納動作の完了成功時に、セグメント・テーブルを更新する。次に、ステップ762において、更新したセグメント・テーブルをカタログ・マネージャに送る。

#### 【0118】

また、一種の冗長情報、例えば、セグメントのコピーを有するファイルを、他の種類の冗長情報、例えば、2つ以上のセグメントの排他的ORに変換することも可能である。例えば、図6に示すプロセスを用いて、データの追加コピーを作成することができる。このプロセスを完了した後、他の形態の冗長データ（セグメントの排他的ORの結果）を削除することができる。同様に、図24に示すプロセスは、格納したデータと共に用いて、冗長情報の排他的ORを作成することも可能である。このような情報を作成した後、図7に示したプロセスを用いてデータの余分なコピーを削除することができる。ファイルが冗長情報を有する形態は、ファイル毎に異なってもよく、更に、例えば、ファイルに関連する優先度に基づくことも可能であり、冗長情報の形態の指示をカタログ・マネージャに格納することも可能である。

#### 【0119】

ネットワークを通じたデータ転送をスケジュールすることにより、そして記憶装置上の負荷を分散し、冗長情報と共にランダムに分散したデータ・セグメントに対するアクセスを選択することによって、このシステムは、多数のアプリケーションと多数の記憶装置との間で双方向に多数のデータ・ストリームを非常にスケーラブルにかつ信頼性高く、効率的に転送することが可能となる。これは、分散型マルチメディア製品にとっては特に有効である。

#### 【0120】

このようなコンピュータ・ネットワークを用いて実現可能なアプリケーションの1つは、多数のストリームを他の外部デジタル効果システムに送りがつ戻し

、ライブ製作において共通に用いる機能である。これらのシステムは複雑でコスト高となる場合がある。殆どのディスクを用いたノンリニア・ビデオ編集システムは、効果リターン・ストリームを同時に記録しながら多数の再生ストリームを維持することができない、ディスク・サブシステムおよびバス・アーキテクチャを有し、オンライン環境で用いる可能性を制限している。このシステムを用いれば、数種類のストリームを効果システムに送ることができ、効果システムは効果データ・ストリームを出力し、多数の記憶装置上に格納する。数種類のストリームは、多数のカメラ・ソース、またはデュアル・ディジタル・ビデオ効果のレイヤとすることができる。

#### 【0121】

また多数の記憶装置を有し、データを1つのクライアントに供給し、他のいずれの記憶装置よりも高い帯域幅を有する高帯域幅データ・ストリームを必要とするクライアントを満足させることも可能である。例えば、20個の記憶装置の各々がスイッチに対して10MB/sリンクを有し、クライアントがスイッチに対して200MB/sリンクを有する場合、クライアントは20個の記憶装置から同時に200MB/sを読み取ることができ、例えば、高品位テレビジョン（HDTV）のデータ・ストリーム転送が可能となる。

#### 【0122】

これまでに概略的に説明した手順を用いると、記憶装置およびクライアントは、ローカル情報を用い、中央のコンフィギュレーション管理または制御なく動作する。記憶装置は、システムを停止する必要なく、動作の間にシステムに追加することができる。記憶装置は単に動作を開始し、クライアントにその可用性を知らせ、次いでアクセス要求に対して応答するプロセスを確立する。この拡張性は、システムの機能および信頼性を補完するものである。

#### 【0123】

以上いくつかの実施形態について説明したが、これまでに述べたことは単なる例示であって限定ではなく、一例として提示したに過ぎないことは、当業者には明白であろう。多数の変更や別の実施形態は当業者の範囲内のことであり、添付した特許請求の範囲およびその均等物に該当するものと見なすこととする。



## 【図面の簡単な説明】

## 【図 1】

図 1 A は、コンピュータ・システムの一例のブロック図である。図 1 B は、図 1 A のシステムの別の実施形態のブロック図である。

## 【図 2】

図 2 A は、図 1 A における記憶装置 4 2 にデータ・セグメントをマッピングするデータ構造を示す。図 2 B は、図 1 B におけるデータ記憶装置 4 2 のセグメントをマッピングするデータ構造を示す。

## 【図 3】

一実施形態において、どのようにしてデータを捕獲し数個の記憶装置間で分散するのかについて記載したフローチャートである。

## 【図 4】

一実施形態において、どのようにして記憶装置はデータ格納要求を処理するのかについて記載したフローチャートである。

## 【図 5】

記憶装置が利用できなくなった場合に、どのようにして故障回復を実行するかについて記載したフローチャートである。

## 【図 6】

追加データ・コピーをどのようにして作成することができるかについて記載したフローチャートである。

## 【図 7】

データ・コピーをどのようにして削除することができるかについて記載したフローチャートである。

## 【図 8】

記憶装置をどのようにしてシステムから除去することができるかについて記載したフローチャートである。

## 【図 9】

どのようにしてデータを保管またはバックアップとしてコピーすることができるかについて記載したフローチャートである。

## 【図 1 0】

カタログ・マネージャに記憶装置の可用性を通知するための、記憶装置上での処理の状態図である。

## 【図 1 1】

カタログ・マネージャが維持することができる記憶装置のリストを示す。

## 【図 1 2】

どのようにしてカタログ・マネージャは記憶装置を監視することができるのかについて示す状態図である。

## 【図 1 3】

メディア・データ・ファイルの同等性を追跡するためのテーブルを示す。

## 【図 1 4】

数個のクリップから成るモーション・ビデオ・シーケンスを表わすリスト構造を示す。

## 【図 1 5】

2つのモーション・ビデオ・データ・ストリーム、および4つの関連するオーディオ・データ・ストリームの再生にクライアントにおいて対応するためのバッファ・メモリの構造を示す。

## 【図 1 6】

クライアントはいかにしてマルチメディア組成を処理し、選択した記憶装置からのデータに対する要求にするのかについて記載したフローチャートである。

## 【図 1 7】

一実施形態において、クライアントはいかにして記憶装置に一次ストレージからバッファにデータ転送を要求するかについて記載したフローチャートである。

## 【図 1 8】

記憶装置はいかにして図 1 7におけるクライアントからの要求に回答するのかについて記載したフローチャートである。

## 【図 1 9】

データに対するディスク・アクセス要求に優先順位を付けるためのディスク・キューの一例、およびネットワーク・データ転送要求に優先順位を付けるための

ネットワーク・キューを示す。

【図20】

一実施形態において、クライアントがどのようにして記憶装置にネットワークを通じたデータ転送を要求するかについて記載したフローチャートである。

【図21】

一実施形態において、記憶装置がどのようにして多数のクライアントからのデータ転送要求を処理するかについて記載するフローチャートである。

【図22】

クライアントがデータを当該クライアントから記憶装置に転送する際に実行するネットワーク・スケジューリング・プロセスの一実施形態について記載するフロー・チャートである。

【図23】

記憶装置がデータをクライアントから当該記憶装置に転送する際に実行するネットワーク・スケジューリング・プロセスの一実施形態について記載するフロー・チャートである。

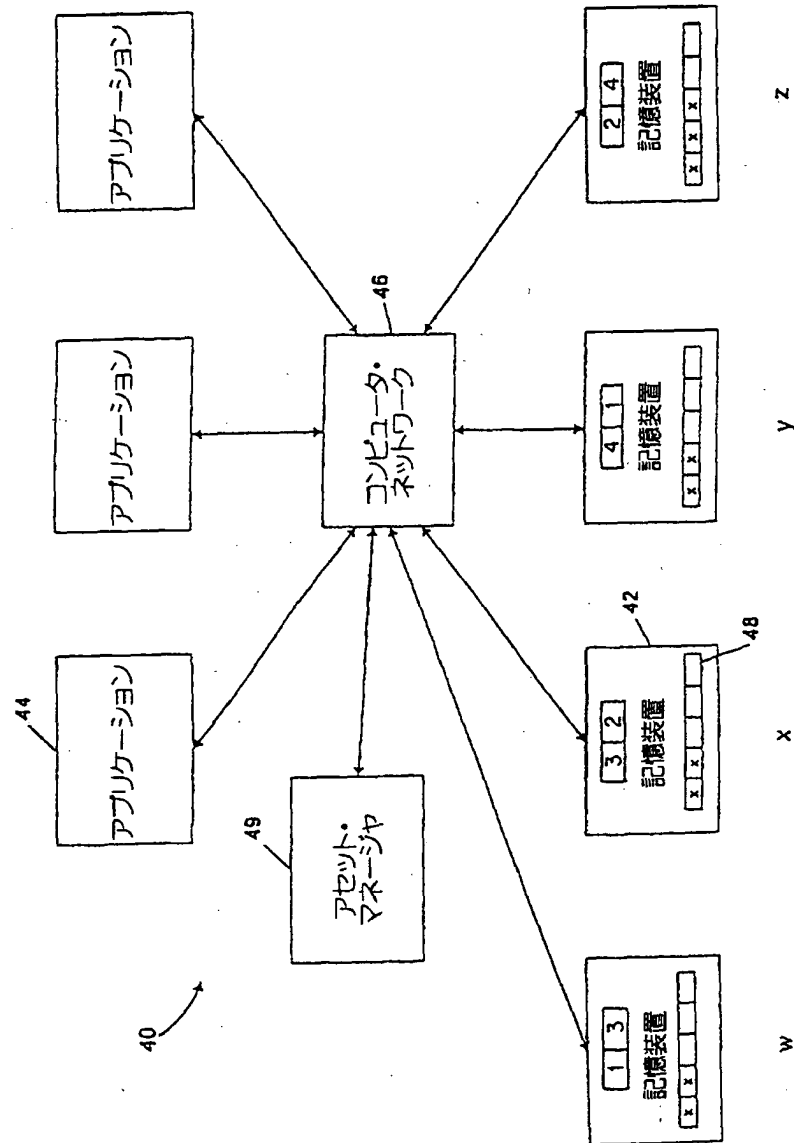
【図24】

別の実施形態において、どのようにしてデータを捕獲し、数個の記憶装置間で分散することができるのかについて記載したフロー・チャートである。

【図25】

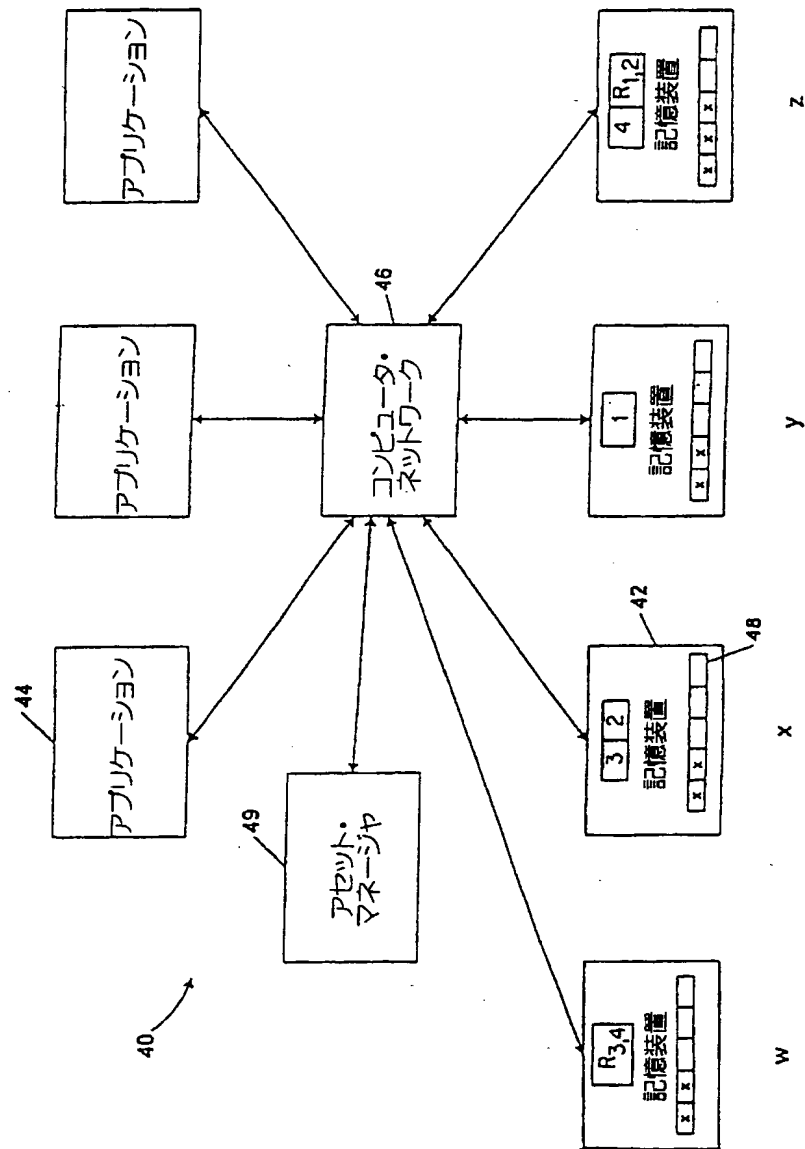
別の実施形態において、記憶装置が使用できなくなった場合に、どのようにして故障回復を実行することができるのかについて記載したフロー・チャートである。

【図1A】



【FIG. 1A】

【図 1 B】



【FIG. 1B】

【図 2 A】

90A

94A

92A

	A	B	.
1	W	Y	...
2	Z	X	...
3	X	W	...
4	Y	Z	...
.	.	.	
.	.	.	
.	.	.	

FIG. 2A

【図 2 B】

90B

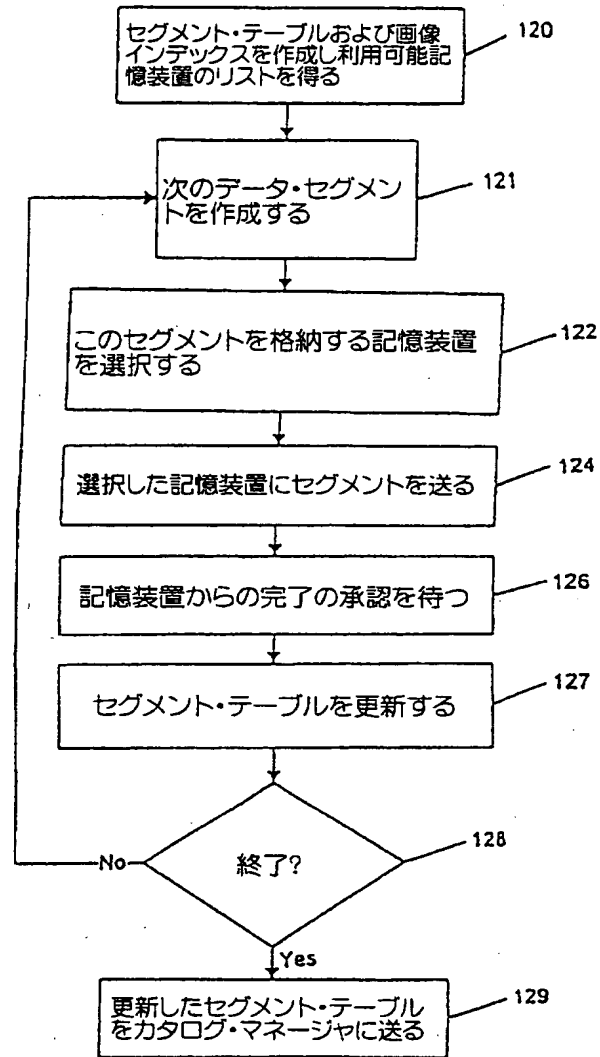
94B    96B

92B

	A	R	.
1	Y	$R_{1,2}$	...
2	X	$R_{1,2}$	...
$R_{1,2}$	Z	$\{1,2\}$	...
3	X	$R_{3,4}$	...
4	Z	$R_{3,4}$	...
$R_{3,4}$	W	$\{1,2\}$	...
.	.	.	...

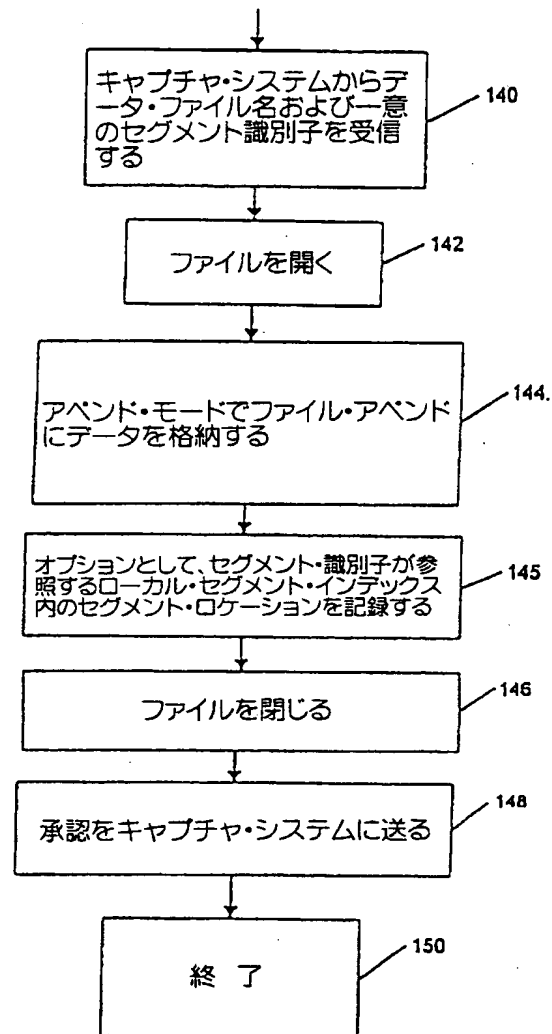
FIG. 2B

【図3】



【 FIG. 3】

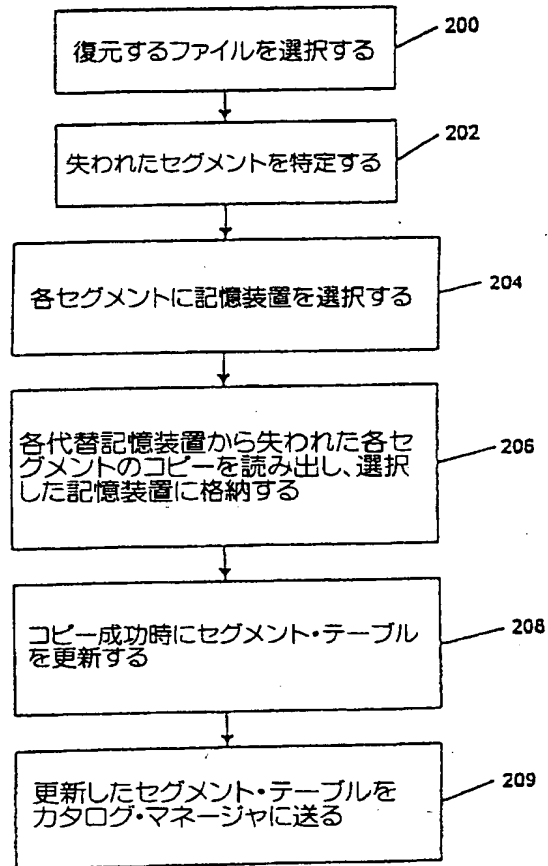
【図4】



【 FIG. 4 】

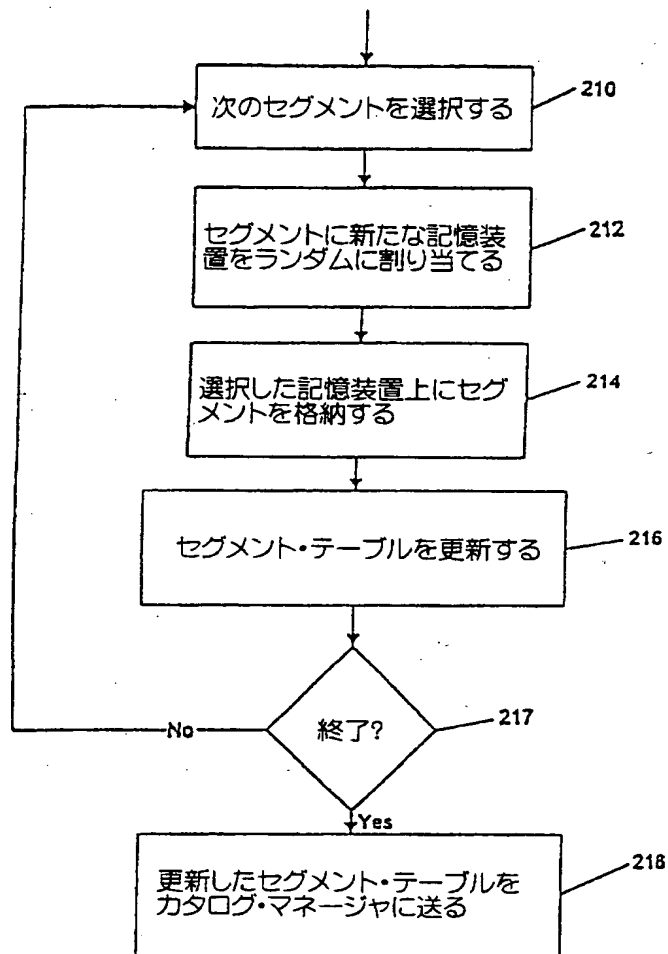


【図 5】



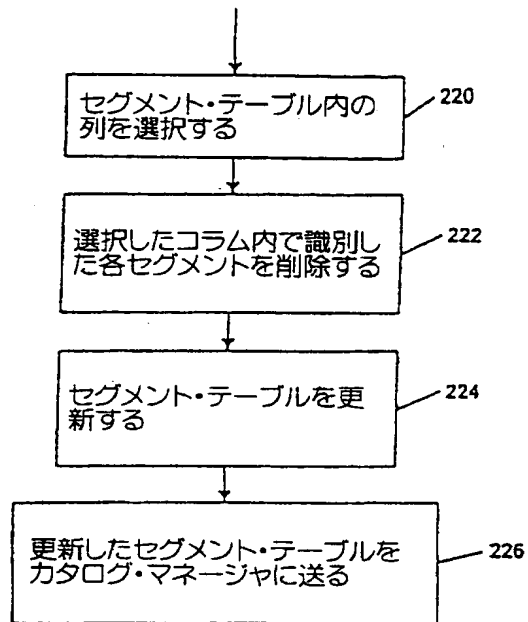
【 FIG. 5】

【図 6】



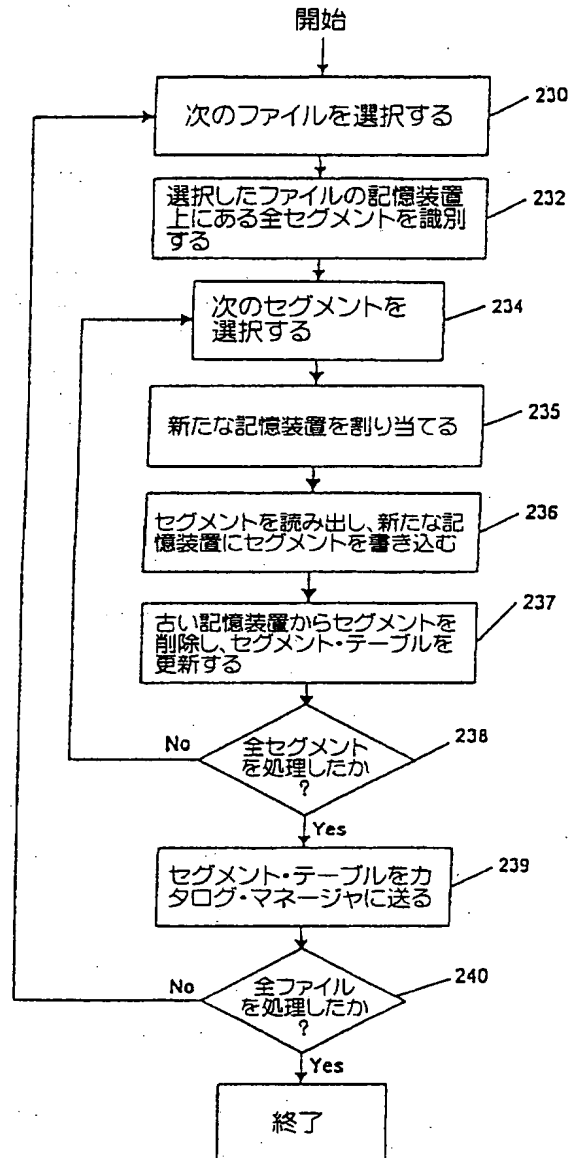
【 FIG. 6 】

【図 7】



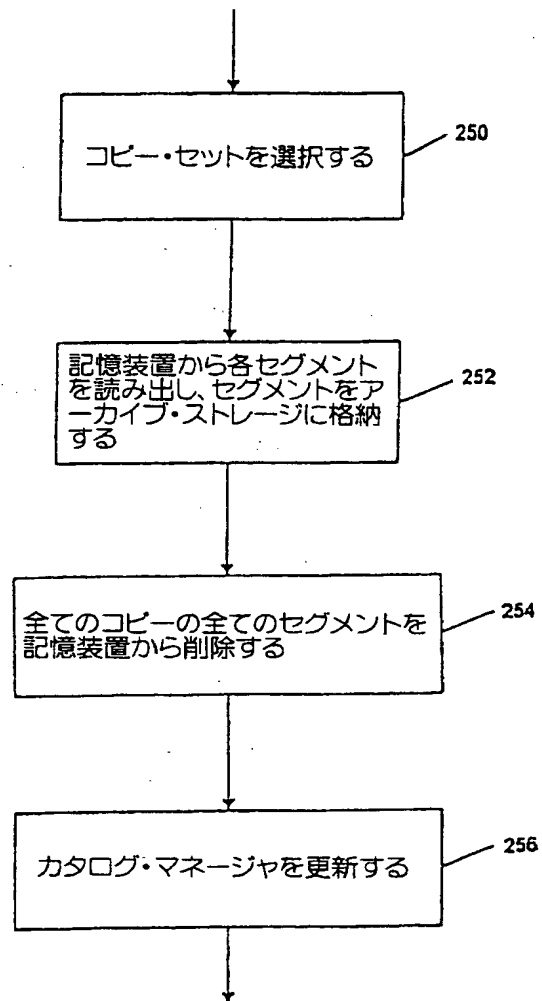
【 FIG. 7】

【図 8】



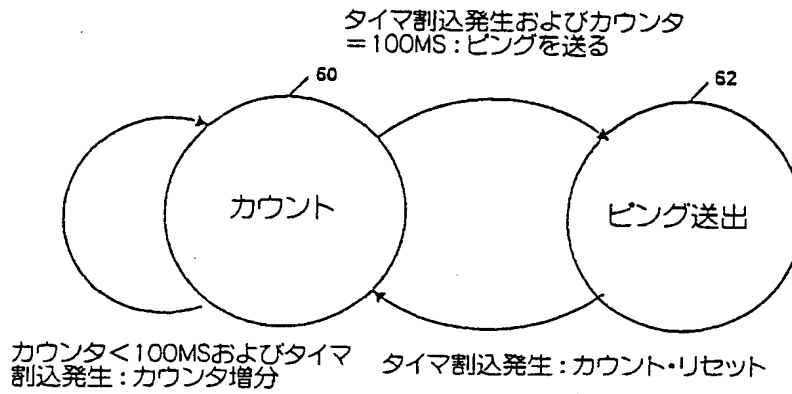
【 FIG. 8 】

【図 9】



【 FIG. 9】

【図10】



【FIG. 10】

【図11】

記憶装置ID

1	帯域幅、メモリ、 容量.....	最後のピング以降のカウント
2		
3		
.		
.		
.		
N		

記憶装置のリスト

70

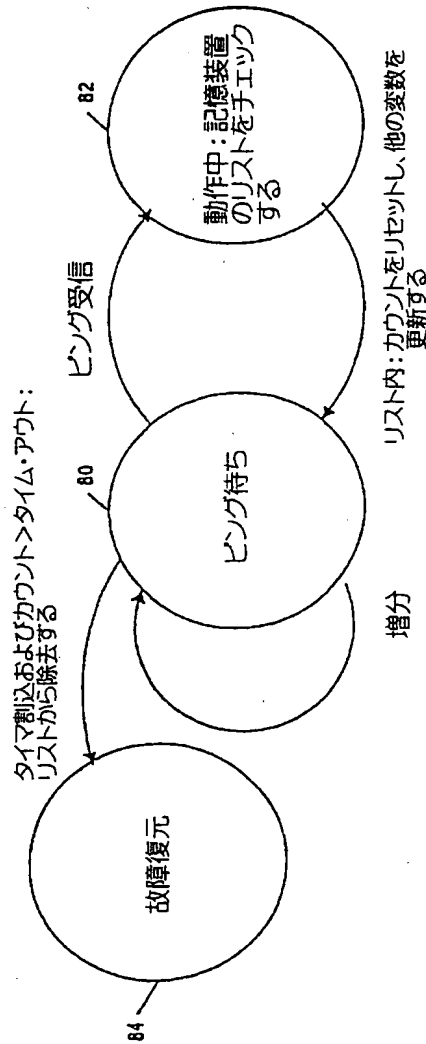
72

74

76

【FIG. 11】

【図 1 2】



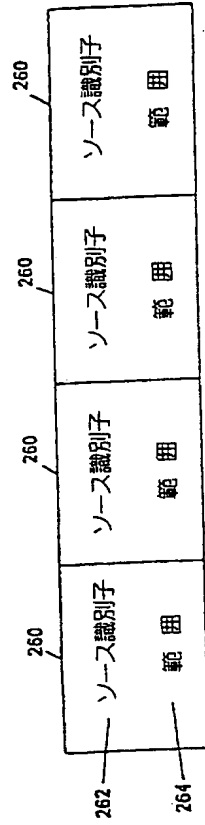
【FIG. 12】

【図 1 3】

100 ソース識別子、範囲	102	104	106	108
ファイル 1		A リスト 1	B リスト 1	
ファイル 2		A リスト 2	B リスト 2	
ファイル 3		A リスト 3	B リスト 3	

【FIG. 13】

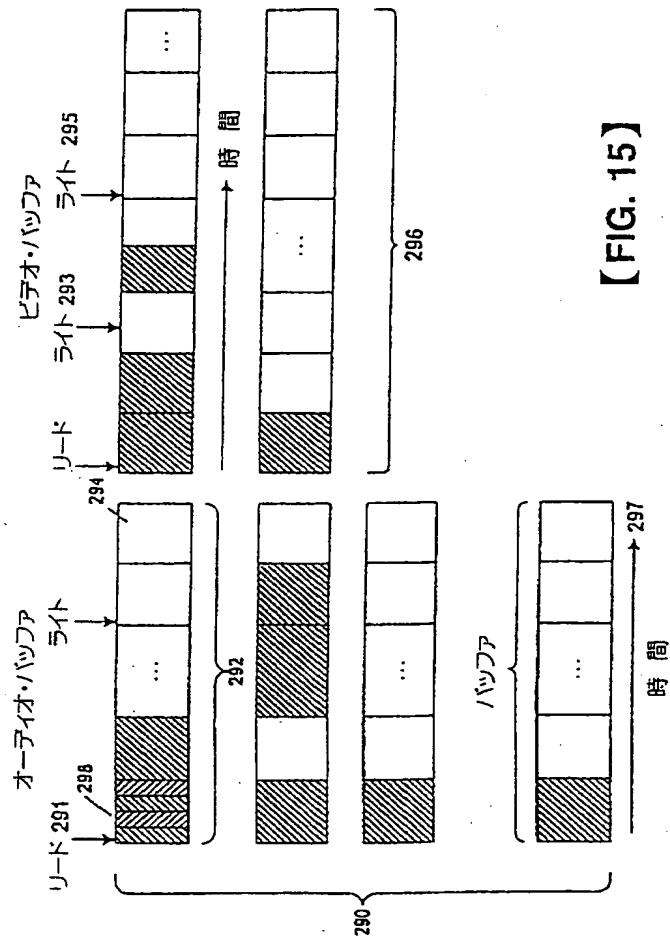
【図 1 4】



【FIG. 14】

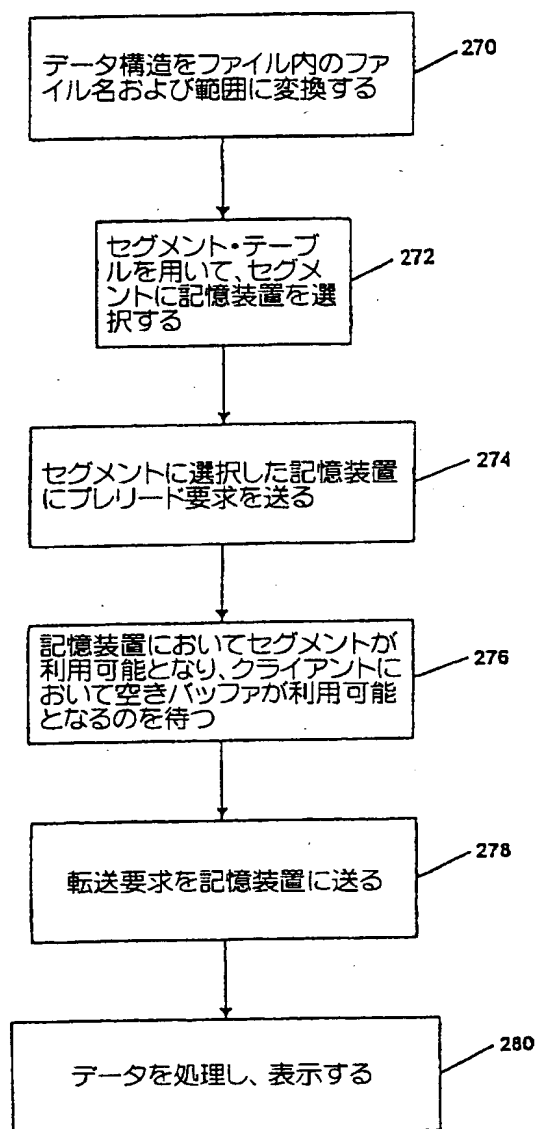


【図 15】



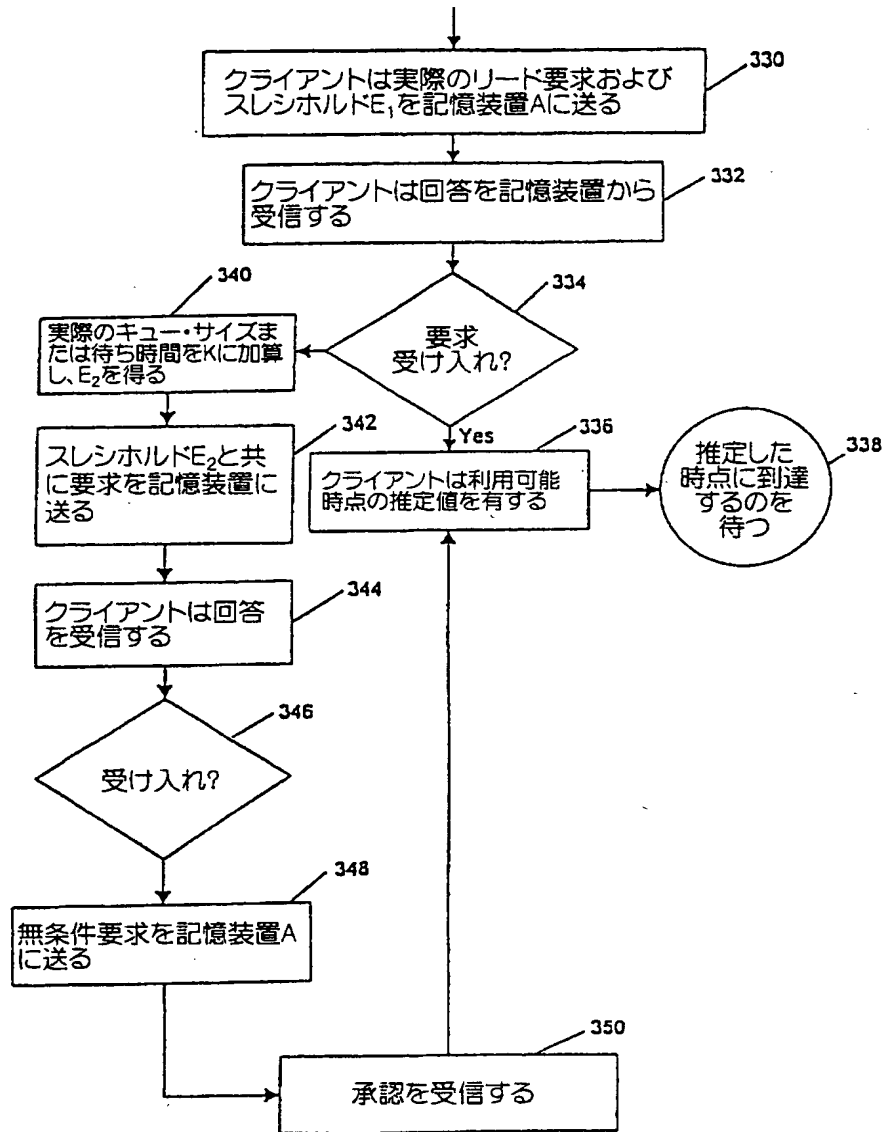
【FIG. 15】

【図 16】



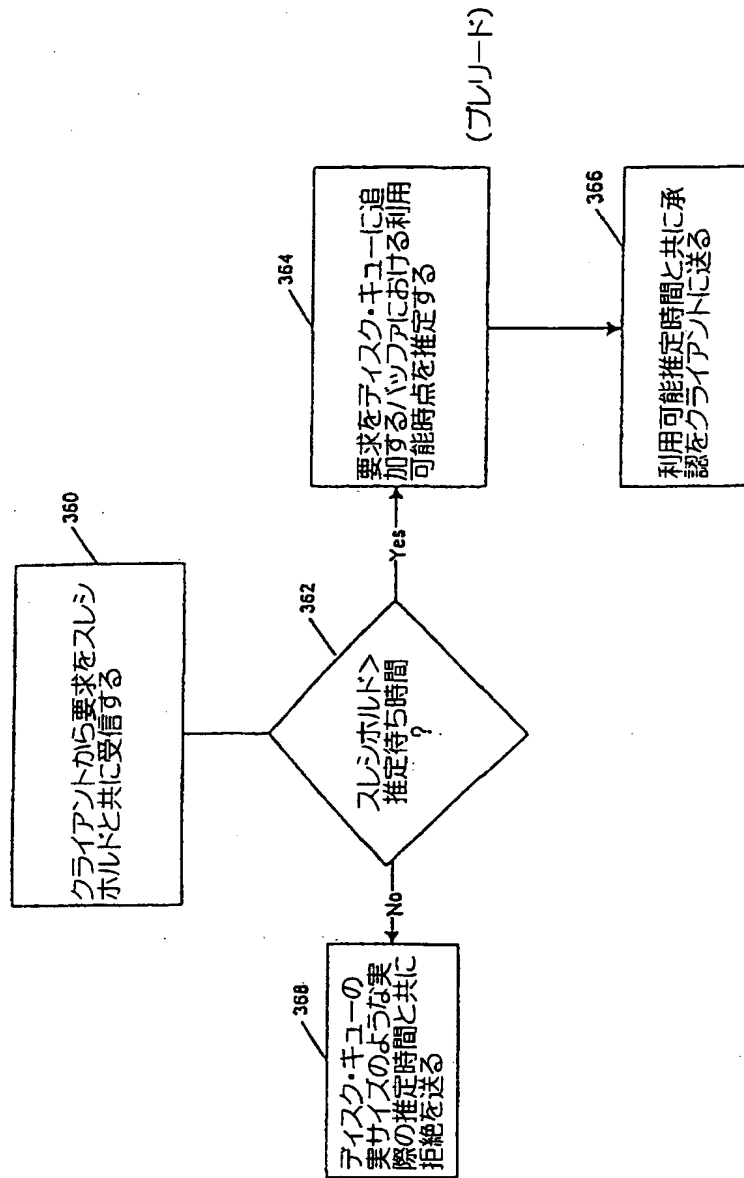
【 FIG. 16】

【図17】



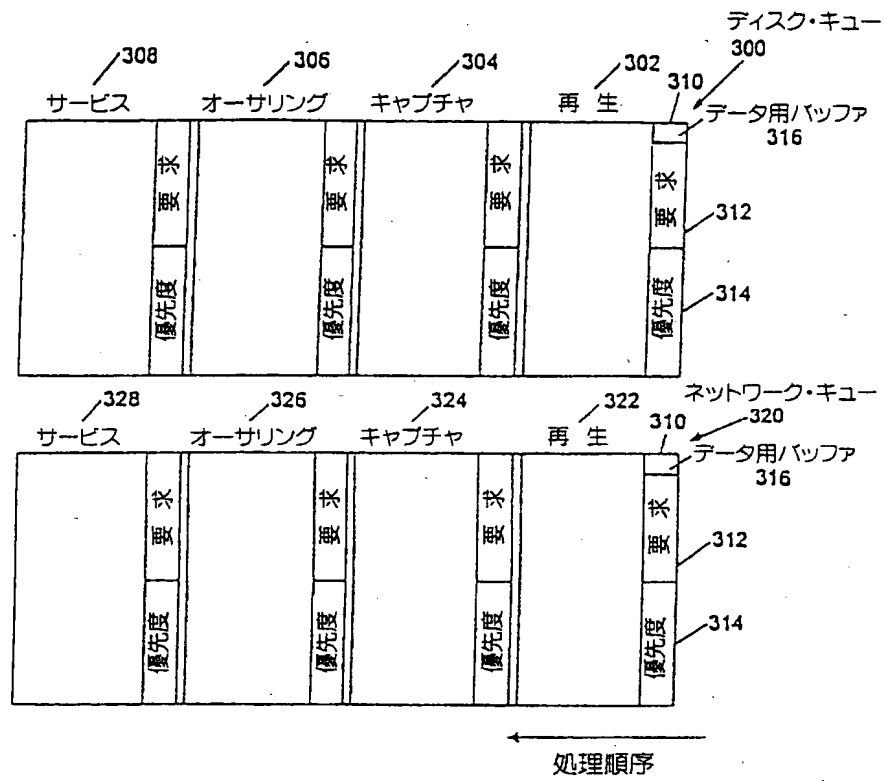
【FIG. 17】

【図18】



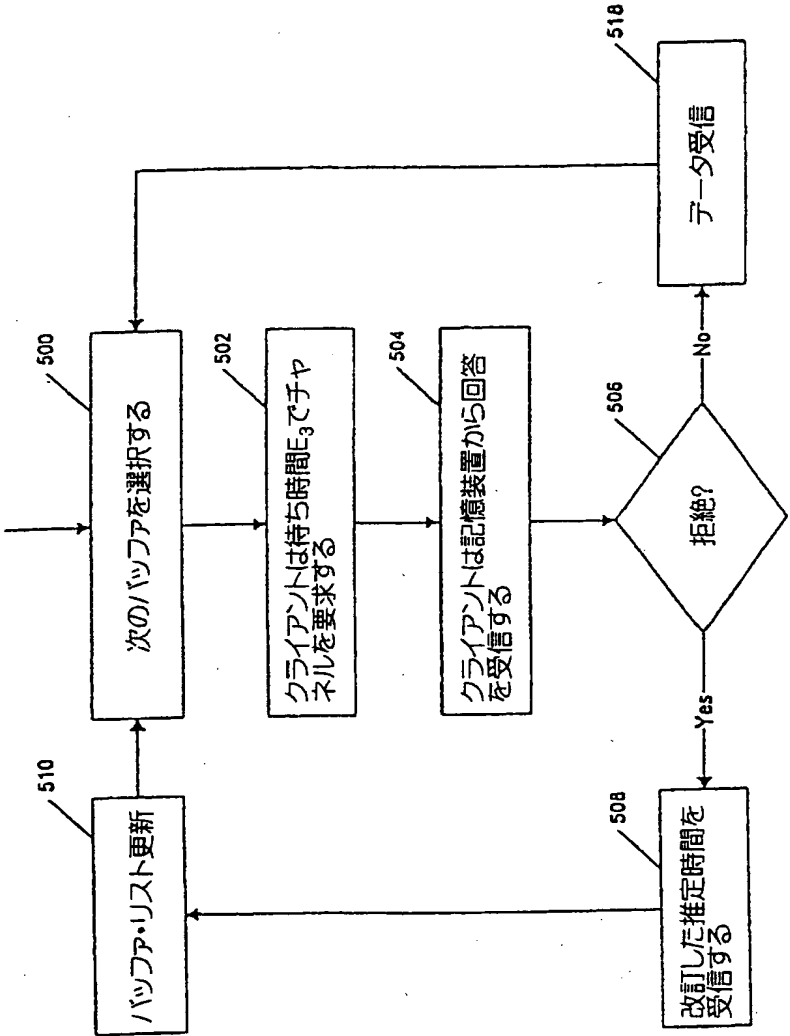
【FIG. 18】

【図19】



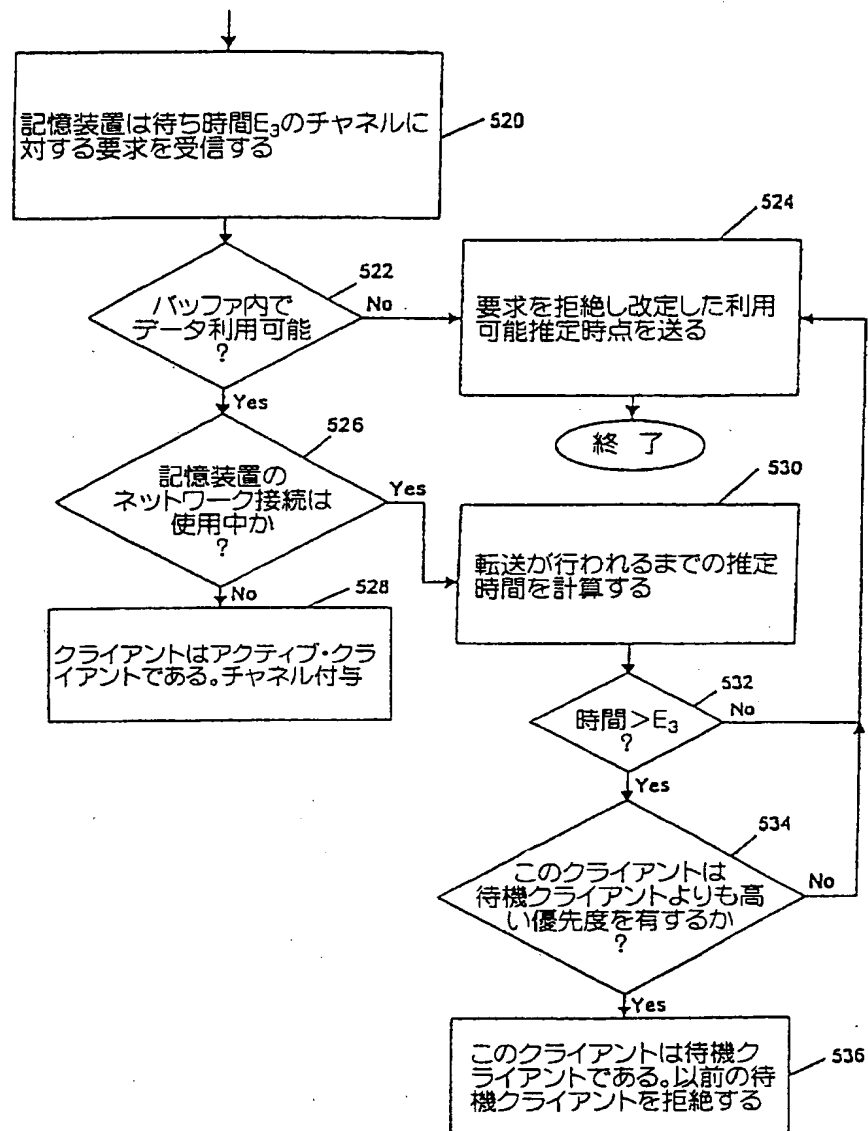
【 FIG. 19 】

【図20】



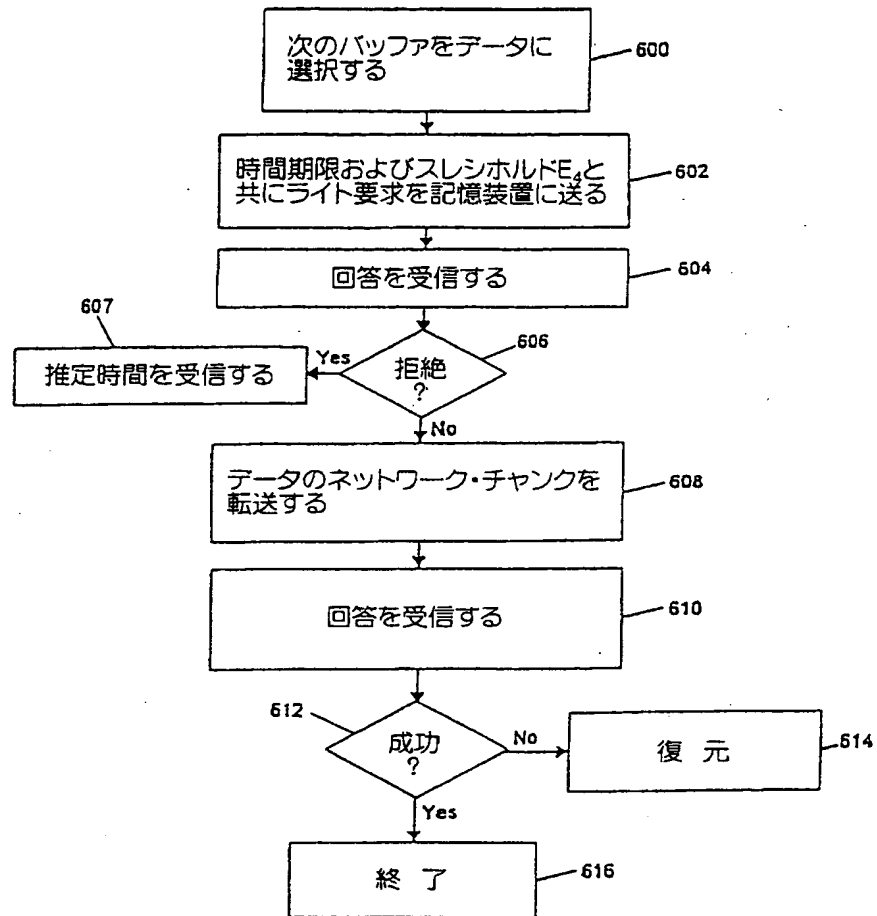
【FIG. 20】

【図 2 1】



【 FIG. 21】

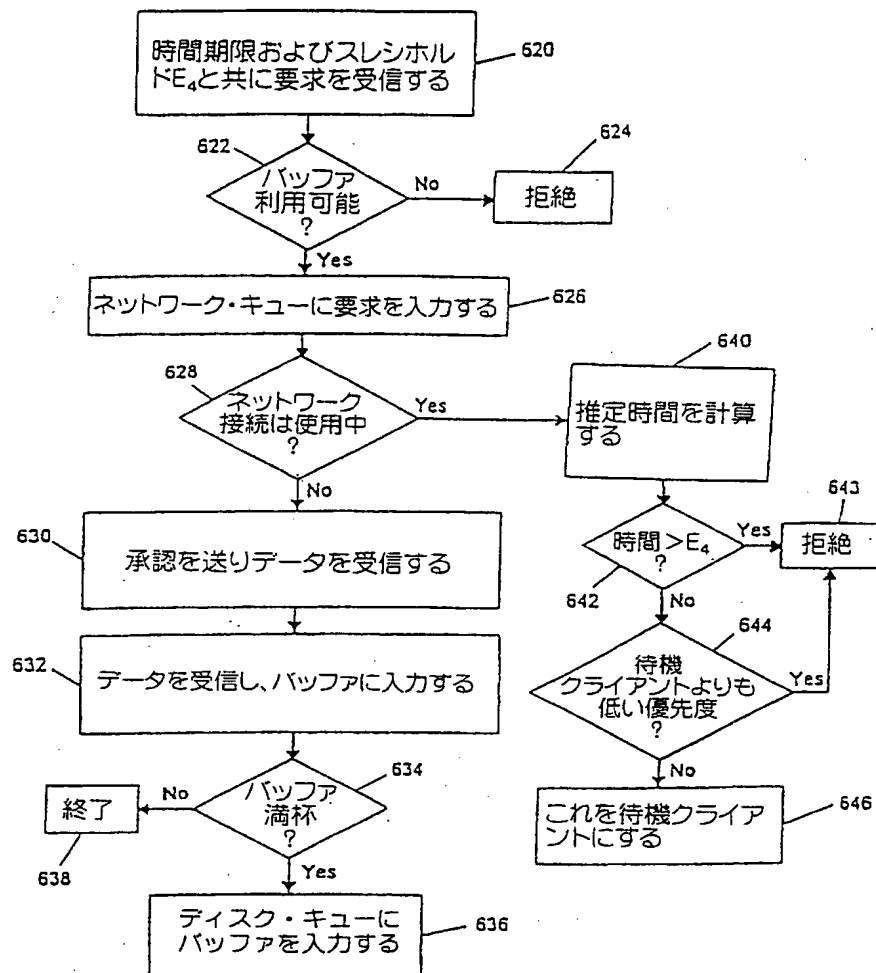
【図 2 2】



【 FIG. 22】

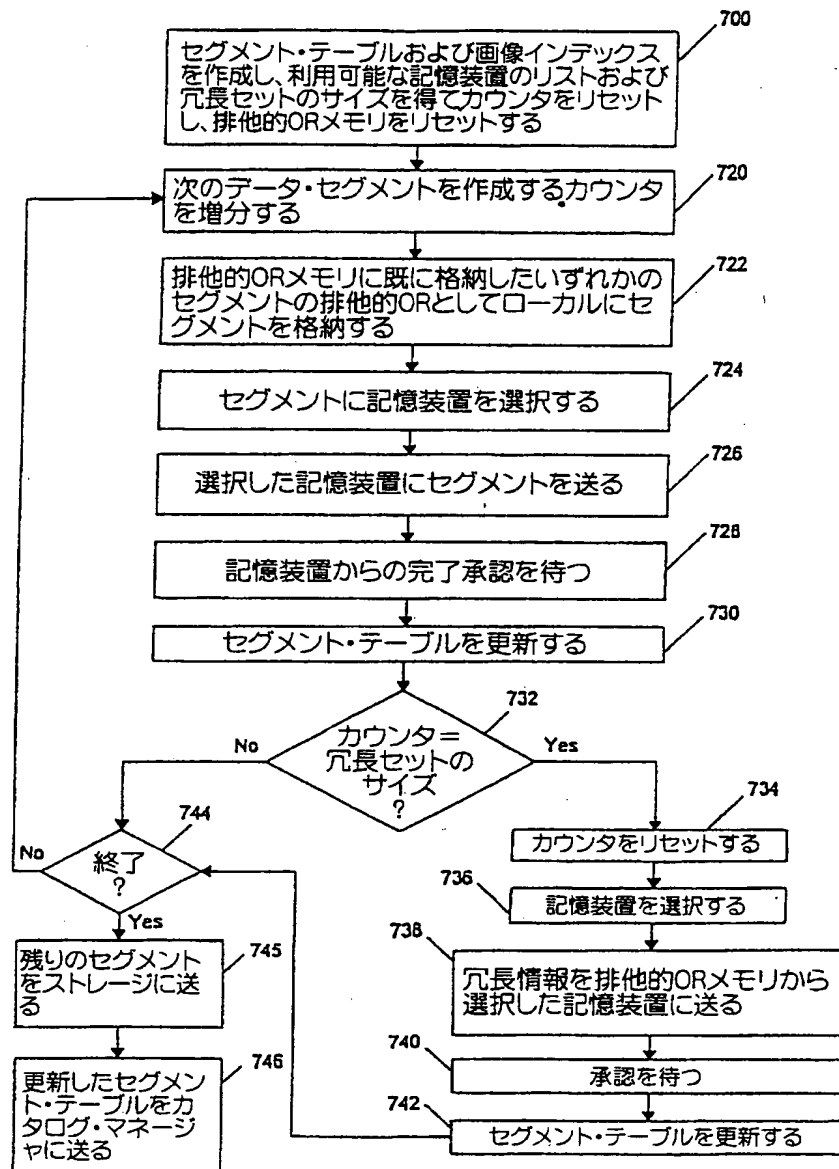


【図23】



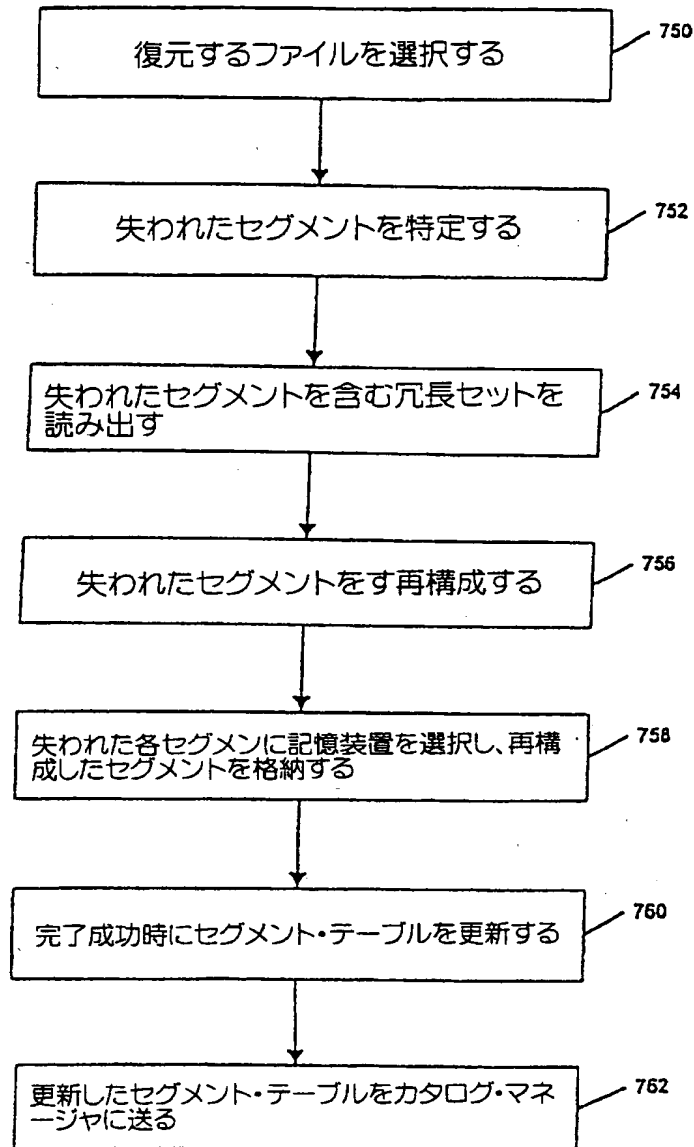
【FIG. 23】

【図24】



【FIG. 24】

【図 25】



【 FIG. 25】

【手続補正書】特許協力条約第34条補正の翻訳文提出書

【提出日】平成12年4月7日(2000.4.7)

【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】特許請求の範囲

【補正方法】変更

【補正内容】

【特許請求の範囲】

【請求項1】 コンピュータ用ファイル・システムであって、前記コンピュータが、当該コンピュータ上で実行するアプリケーションからの要求に应答してコンピュータ・ネットワークを通じて独立したリモート記憶装置にアクセスし、前記記憶装置上に格納してあるデータを読み出すことを可能とし、ファイルは、データ・セグメントと各セグメント毎に対応する冗長情報とを含み、各ファイル毎に、各データ・セグメントを前記記憶装置のランダムまたは疑似ランダムに選択した1つに格納し、各データ・セグメント毎に、前記記憶装置のランダムまたは疑似ランダムに選択した1つに前記対応する冗長情報を格納し、セグメントに対する前記冗長情報は、前記セグメントの少なくとも1つのコピーであり、

前記要求に应答してデータを読み出し、前記選択したデータの各セグメント毎に、前記セグメントを格納する記憶装置の1つを選択する手段と、

前記要求したデータの各セグメントを、当該セグメントに対して選択した記憶装置から、読み出す手段と、

各セグメント毎に、前記選択した記憶装置から読み出した選択データのセグメントを直列化する手段と、

前記特定した記憶装置からデータを受信したとき、前記アプリケーションに前記データを供給する手段と、

を備えるファイル・システム。

【請求項2】 前記記憶装置から1つを選択する前記手段は、前記複数の記憶装置上の要求の負荷が実質的に均衡化するように、前記記憶装置を選択する請求項1記載のファイル・システム。

【請求項3】 前記選択手段は、前記セグメントのためのどの記憶装置が、前記要求に応じるための推定時間が最も短いかに基づいて、前記セグメントのための前記記憶装置を選択する請求項2記載のファイル・システム。

【請求項4】 前記選択手段が、  
前記ファイル・システムにおいて、  
前記記憶装置の1つからデータを要求し、推定時間を示す手段と、  
前記第1記憶装置が前記要求を拒絶した場合、前記記憶装置の他のものからデータを要求し、推定時間を示す手段と、  
前記第2記憶装置が前記要求を拒絶した場合、前記第1記憶装置から前記データを要求する手段と、  
を含み、  
各記憶装置において、  
前記要求が前記推定時間以内に前記記憶装置によって対応することができない場合、データの要求を拒絶する手段と、  
前記要求が前記推定時間以内に前記記憶装置によって対応することができる場合、データの要求を受け入れる手段と、  
を含む請求項3記載のファイル・システム。

【請求項5】 各セグメントを読み出す前記手段は、前記選択した記憶装置からの前記データの転送をスケジュールし、前記記憶装置が効率的にデータ転送を行なうようにする手段を備える請求項2記載のファイル・システム。

【請求項6】 転送をスケジュールする前記手段は、  
前記ファイル・システムにおいて、  
前記選択した記憶装置から前記データの転送を要求し、待ち時間を示す手段と、  
前記選択した記憶装置が前記データを転送する前記要求を拒絶した場合、別の記憶装置から前記データを要求する手段と、  
を含み、  
前記記憶装置において、

前記データが前記指示した待ち時間までに前記記憶装置から転送することができない場合、データ転送要求を拒絶する手段と、

前記選択した記憶装置が前記待ち時間以内に前記データを転送することができる場合、前記データを転送する手段と、

を含む請求項5記載のファイル・システム。

【請求項7】 各セグメントを読み出す前記手段は、前記選択した記憶装置からの前記データの転送をスケジュールし、前記記憶装置が効率的にデータ転送を行なうようにする手段を備える請求項1記載のファイル・システム。

【請求項8】 各セグメントを読み出す前記手段は、前記選択した記憶装置からの前記データの転送をスケジュールし、前記記憶装置が効率的にデータ転送を行なうようにする手段を備える請求項1記載のファイル・システム。

【請求項9】 転送をスケジュールする前記手段は、

前記ファイル・システムにおいて、

前記選択した記憶装置から前記データの転送を要求し、待ち時間を示す手段と

前記選択した記憶装置が前記データを転送する前記要求を拒絶した場合、別の記憶装置から前記データを要求する手段と  
を含み、

前記記憶装置において、

前記データが前記指示した待ち時間までに前記記憶装置から転送することができない場合、データ転送要求を拒絶する手段と、

前記選択した記憶装置が前記待ち時間以内に前記データを転送することができる場合、前記データを転送する手段と、

を含む請求項8記載のファイル・システム。

【請求項10】 各セグメントを前記記憶装置の異なるものに格納する請求項1記載のファイル・システム。

【請求項11】 各セグメントの各コピーを、前記記憶装置の相対的仕様の関数として定義した確率分布にしたがって、前記複数の記憶装置の1つに割り当てる請求項10記載のファイル・システム。

【請求項12】 更に、コンピュータ読み取り可能ロジックを格納し、データ・セグメントの指示を用いてコンピュータによるアクセスが可能なセグメント・テーブルを定義するコンピュータ読み取り可能媒体を備え、前記セグメントおよび対応する冗長情報を格納した前記複数の記憶装置から、該記憶装置の指示を検索する請求項1記載の分散データ記憶システム。

【請求項13】 前記複数の記憶装置が、  
前記コンピュータ・ネットワークに接続した第1記憶装置と、  
前記コンピュータ・ネットワークに接続した第2記憶装置と、  
前記コンピュータ・ネットワークに接続した第3記憶装置と、  
を備える請求項1記載のファイル・システム。

【請求項14】 コンピュータ用ファイル・システムであって、前記コンピュータが、当該コンピュータ上で実行するアプリケーションからの要求に応答してコンピュータ・ネットワークを通じて独立したリモート記憶装置にアクセスし、前記記憶装置上にデータを格納することを可能とするファイル・システムにおいて、

前記データを格納する前記要求に応答して、前記データを複数のセグメントに分割する手段と、

各セグメント毎に記憶装置をランダムまたは疑似ランダムに選択し、前記選択した記憶装置に前記データを格納するように要求する手段と、

各セグメント毎に、当該セグメントに対応する冗長情報のための記憶装置をランダムまたは疑似ランダムに選択し、前記選択した記憶装置に前記冗長情報を格納するように要求する手段であって、前記冗長情報が前記セグメントの少なくとも1つのコピーである、手段と、

前記データを格納したか否かについて前記アプリケーションに確認する手段と、  
を備えるファイル・システム。

【請求項15】 記選択する手段は、前記記憶装置の部分集合を選択する手段と、前記選択した部分集合内にある前記記憶装置の中から少なくとも2つの前記記憶装置を選択する手段とを含む請求項14記載のファイル・システム。

【請求項16】 各セグメントの各コピーは、前記記憶装置の相対的仕様の関数として定義した、確率分布にしたがって、前記複数の記憶装置の1つに割り当てられる請求項14記載のファイル・システム。

【請求項17】 更に、コンピュータ読み取り可能ロジックを格納し、データ・セグメントの指示を用いてコンピュータによるアクセスが可能なセグメント・テーブルを定義するコンピュータ読み取り可能媒体を備え、前記セグメントおよび対応する冗長情報を格納した前記複数の記憶装置から、該記憶装置の指示を検索する請求項14記載のファイル・システム。

【請求項18】 前記複数の記憶装置が、  
前記コンピュータ・ネットワークに接続した第1記憶装置と、  
前記コンピュータ・ネットワークに接続した第2記憶装置と、  
前記コンピュータ・ネットワークに接続した第3記憶装置と、  
を備える請求項14記載のファイル・システム。



【手續補正書】

【提出日】平成13年3月5日(2001.3.5)

【手續補正1】

【補正対象書類名】図面

【補正対象項目名】図1

【補正方法】変更

【補正内容】

【図1】

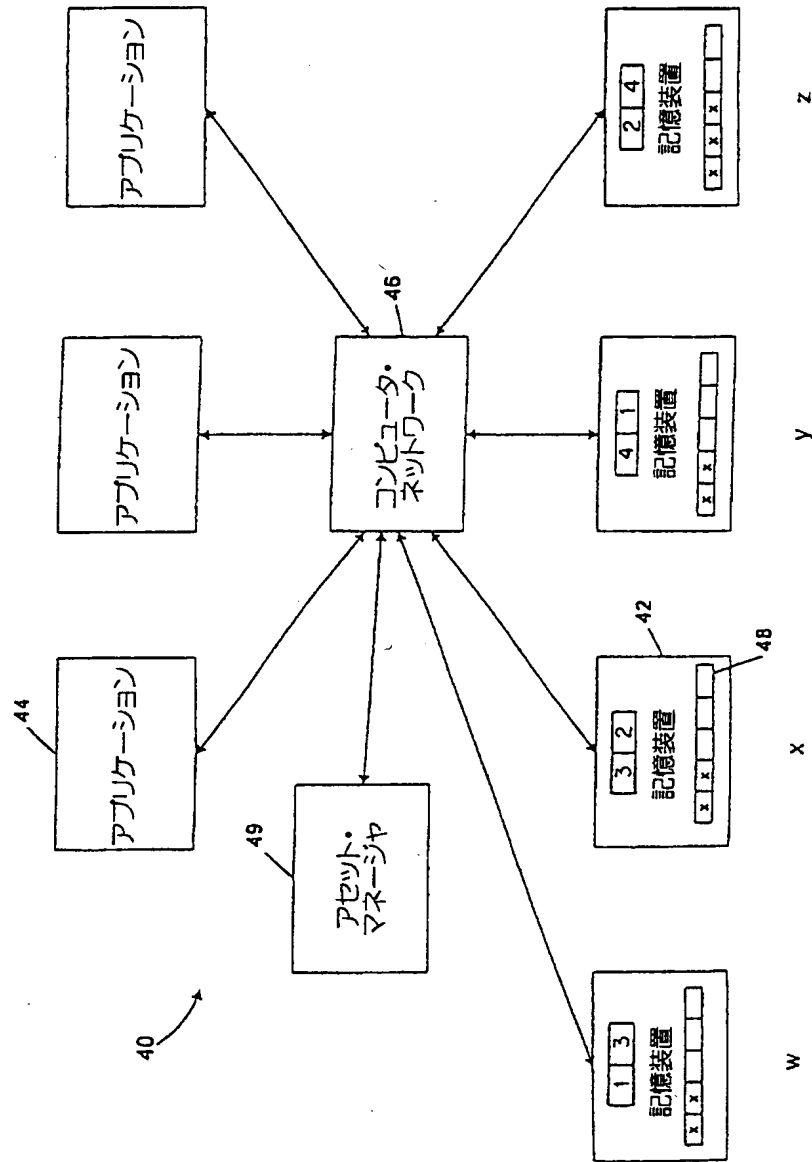


FIG. 1A

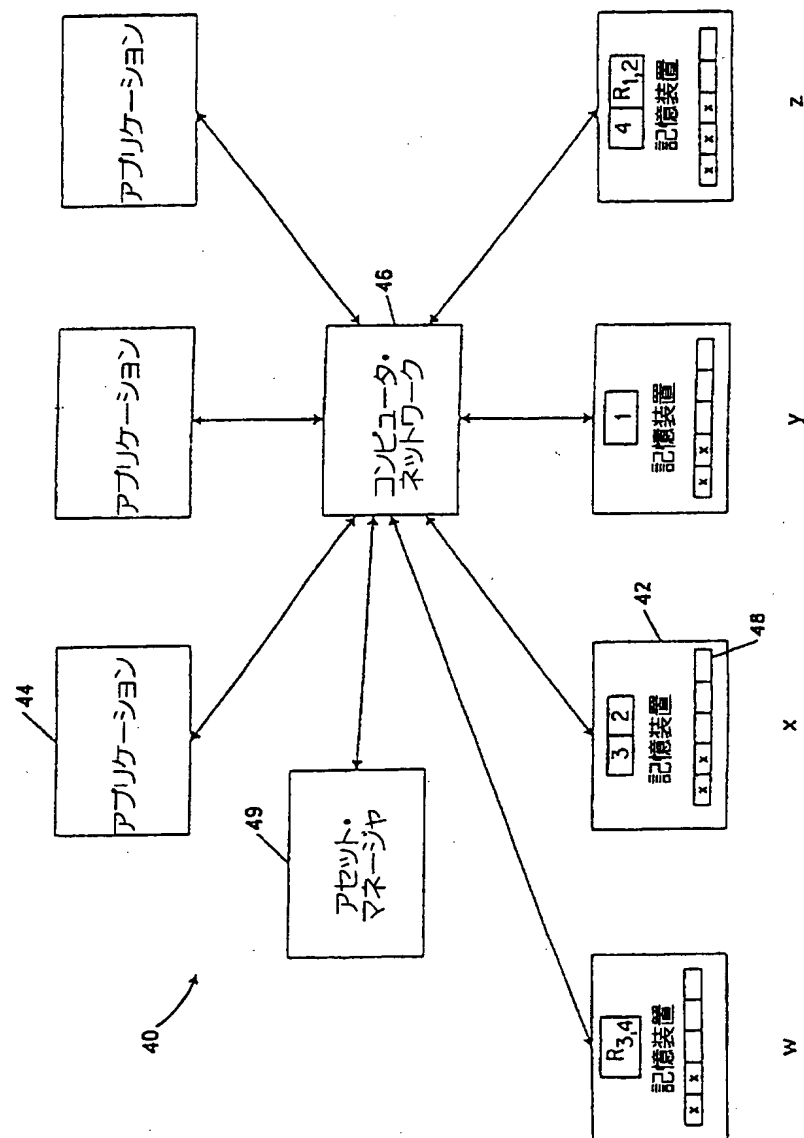


FIG. 1B

【手続補正2】

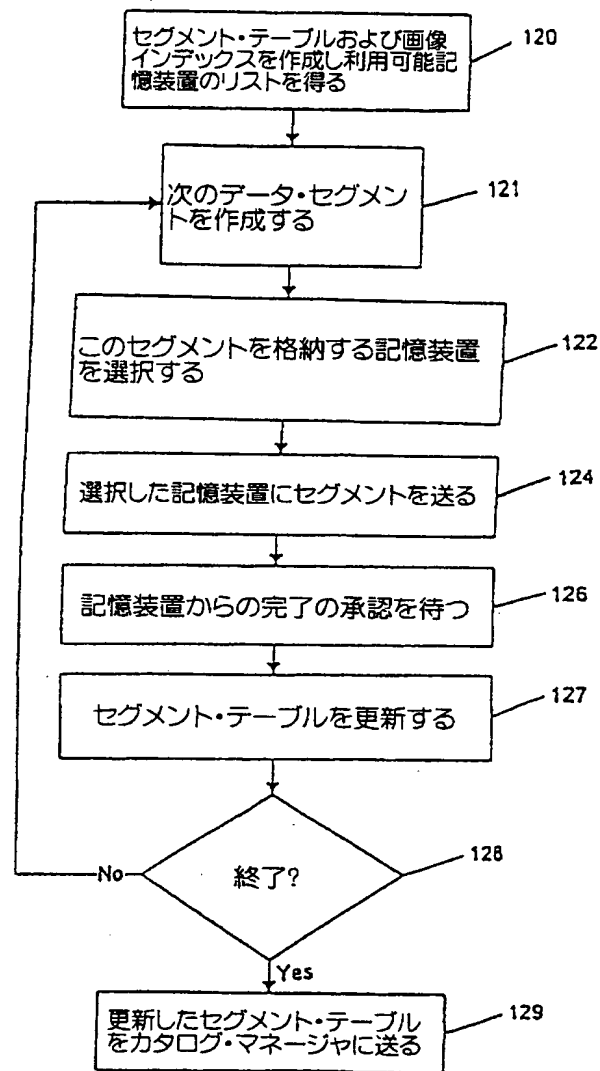
【補正対象書類名】図面

【補正対象項目名】図3

【補正方法】変更

【補正内容】

【図3】



【手続補正3】

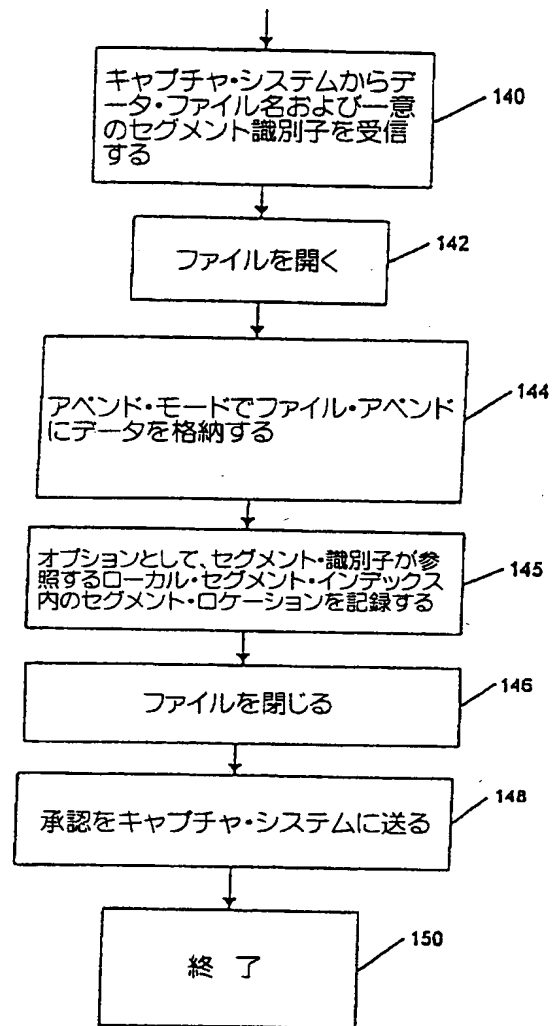
【補正対象書類名】図面

【補正対象項目名】図4

【補正方法】変更

【補正内容】

【図4】



【手続補正4】

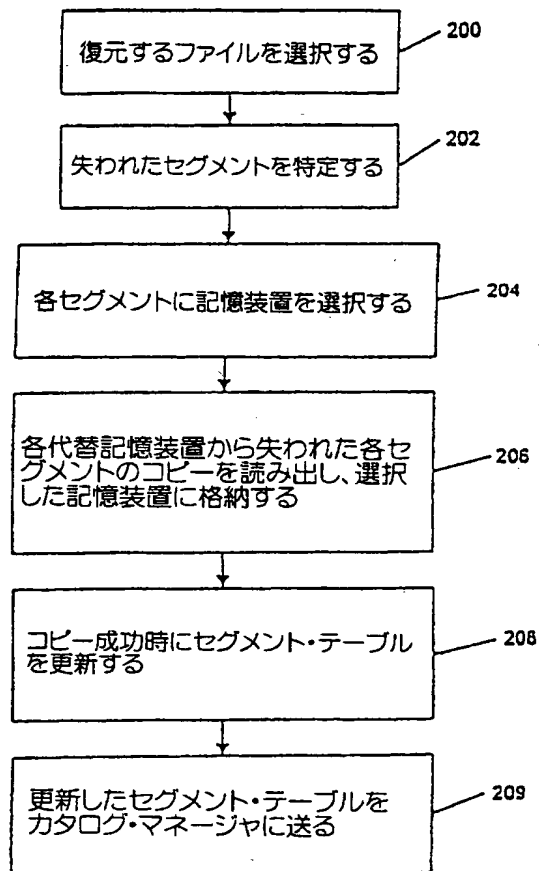
【補正対象書類名】図面

【補正対象項目名】図5

【補正方法】変更

【補正内容】

【図5】



【手続補正5】

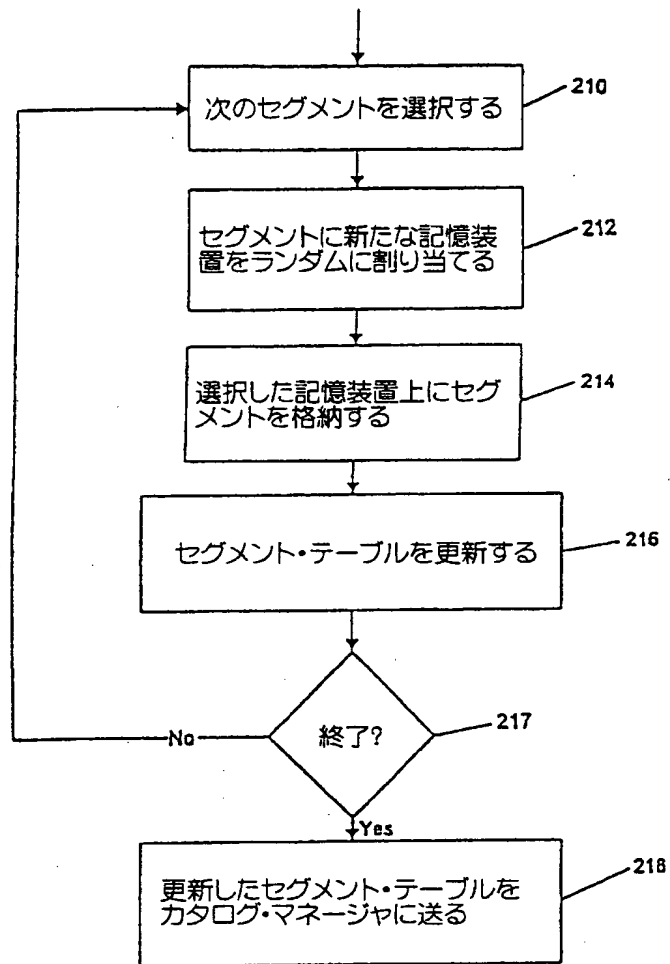
【補正対象書類名】図面

【補正対象項目名】図6

【補正方法】変更

【補正内容】

【図6】



【手続補正6】

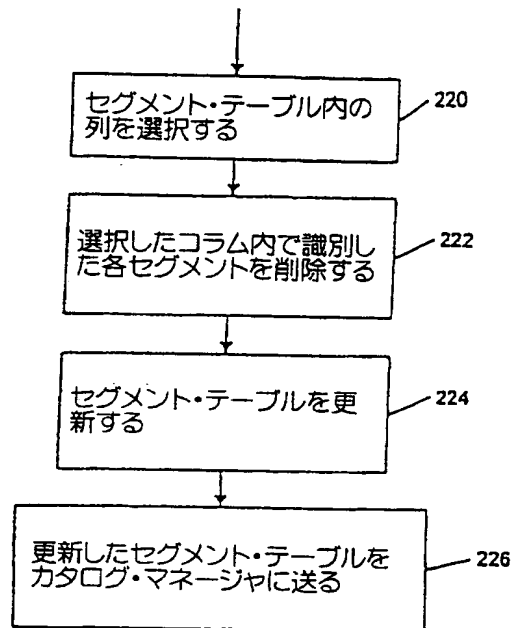
【補正対象書類名】図面

【補正対象項目名】図7

【補正方法】変更

【補正内容】

【図7】



【手続補正7】

【補正対象書類名】図面

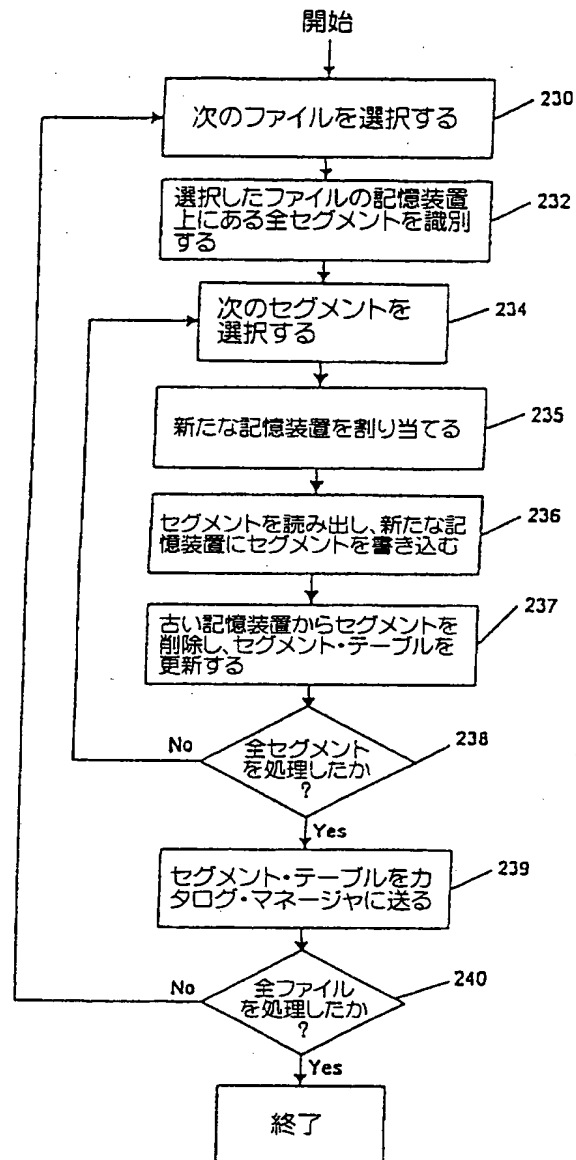
【補正対象項目名】図8

【補正方法】変更

【補正内容】



【図 8】



【手続補正 8】

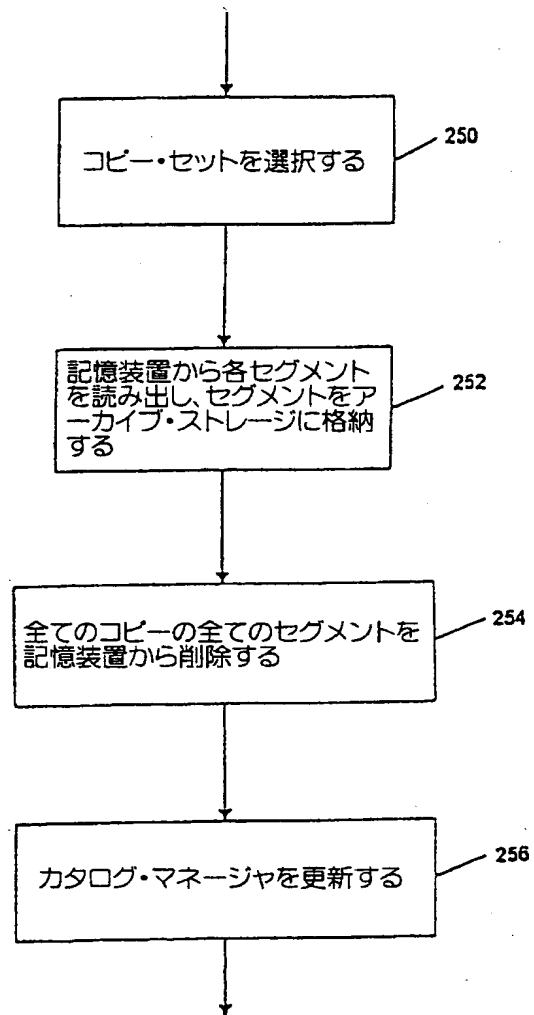
【補正対象書類名】図面

【補正対象項目名】図 9

【補正方法】変更

【補正内容】

【図9】



【手続補正9】

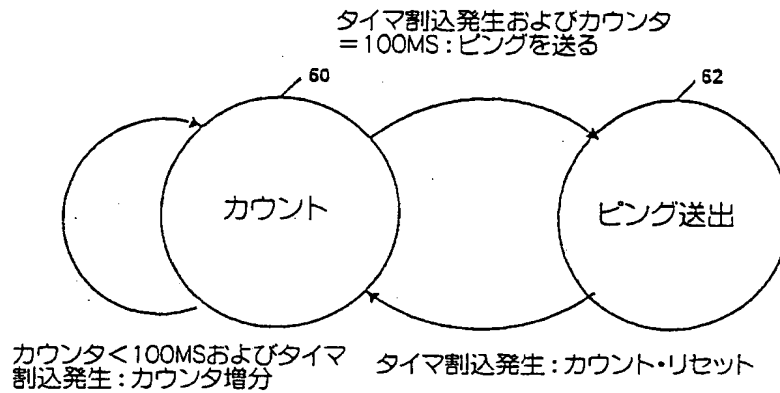
【補正対象書類名】 図面

【補正対象項目名】 図10

【補正方法】 変更

【補正内容】

【図10】



【手続補正10】

【補正対象書類名】図面

【補正対象項目名】図11

【補正方法】変更

【補正内容】

【図11】

記憶装置ID

1	帯域幅、メモリ、容量……	最後のピング以降のカウント
2		
3		
.		
.		
.		
N		

記憶装置のリスト

70

72 74 76

【手続補正11】

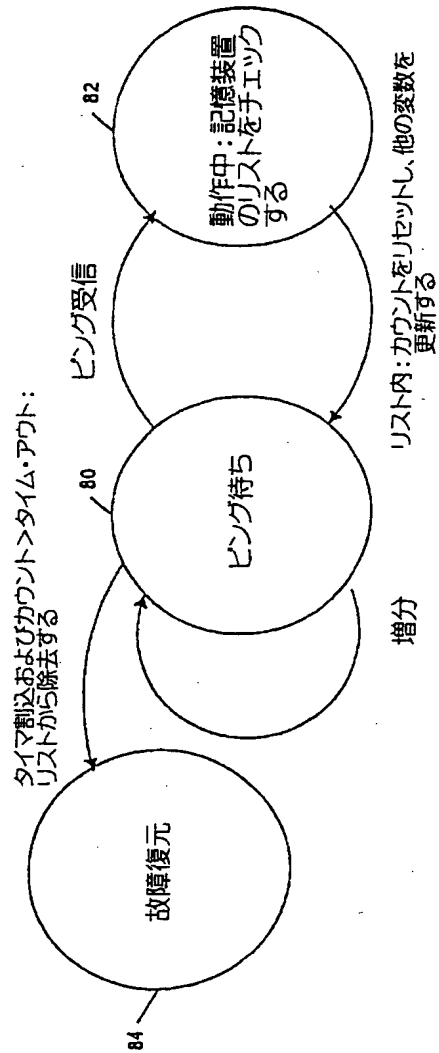
【補正対象書類名】図面

【補正対象項目名】図12

【補正方法】変更

【補正内容】

【図12】



【手続補正12】

【補正対象書類名】図面

【補正対象項目名】図13

【補正方法】変更

【補正内容】

【図13】

100 102 ソース識別子、範囲	104	106	108
	ファイル1	Aリスト1	Bリスト1
	ファイル2	Aリスト2	Bリスト2
	ファイル3	Aリスト3	Bリスト3

【手続補正13】

【補正対象書類名】図面

【補正対象項目名】図14

【補正方法】変更

【補正内容】

【図14】

260	ソース識別子 範囲
260	ソース識別子 範囲
260	ソース識別子 範囲
260	ソース識別子 範囲
262	—
264	—

【手続補正14】

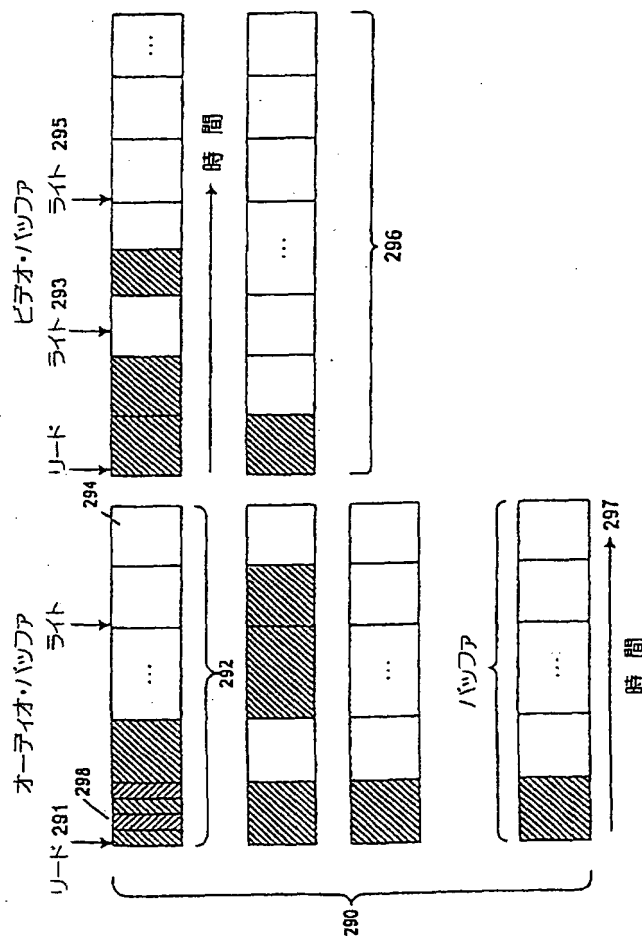
【補正対象書類名】図面

【補正対象項目名】図15

【補正方法】変更

【補正内容】

【図15】



【手続補正15】

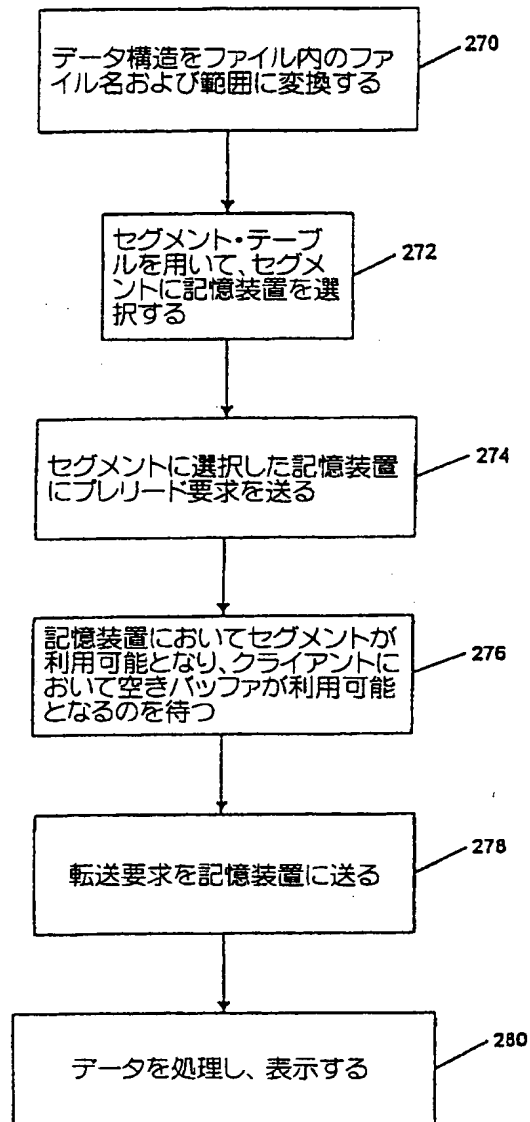
【補正対象書類名】図面

【補正対象項目名】図16

【補正方法】変更

【補正内容】

【図16】



【手続補正16】

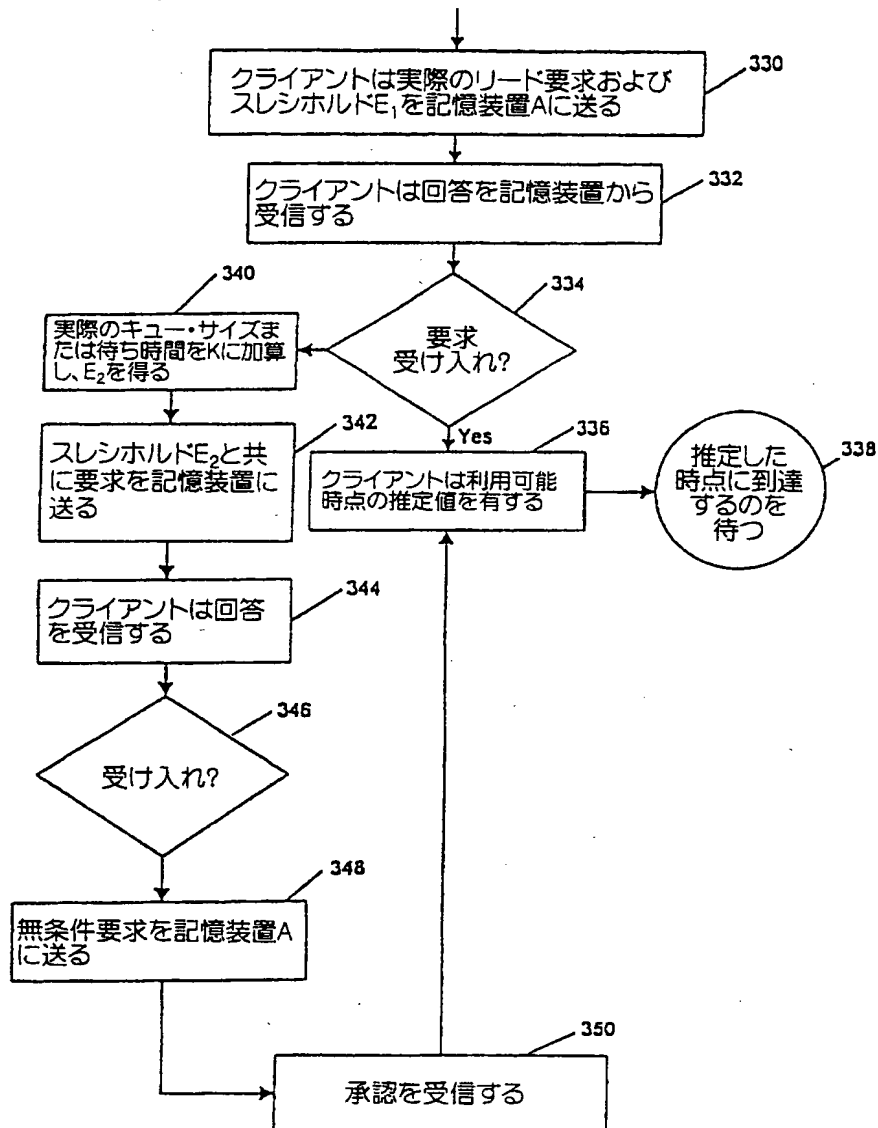
【補正対象書類名】図面

【補正対象項目名】図17

【補正方法】変更

【補正内容】

【図17】



【手続補正17】

【補正対象書類名】図面

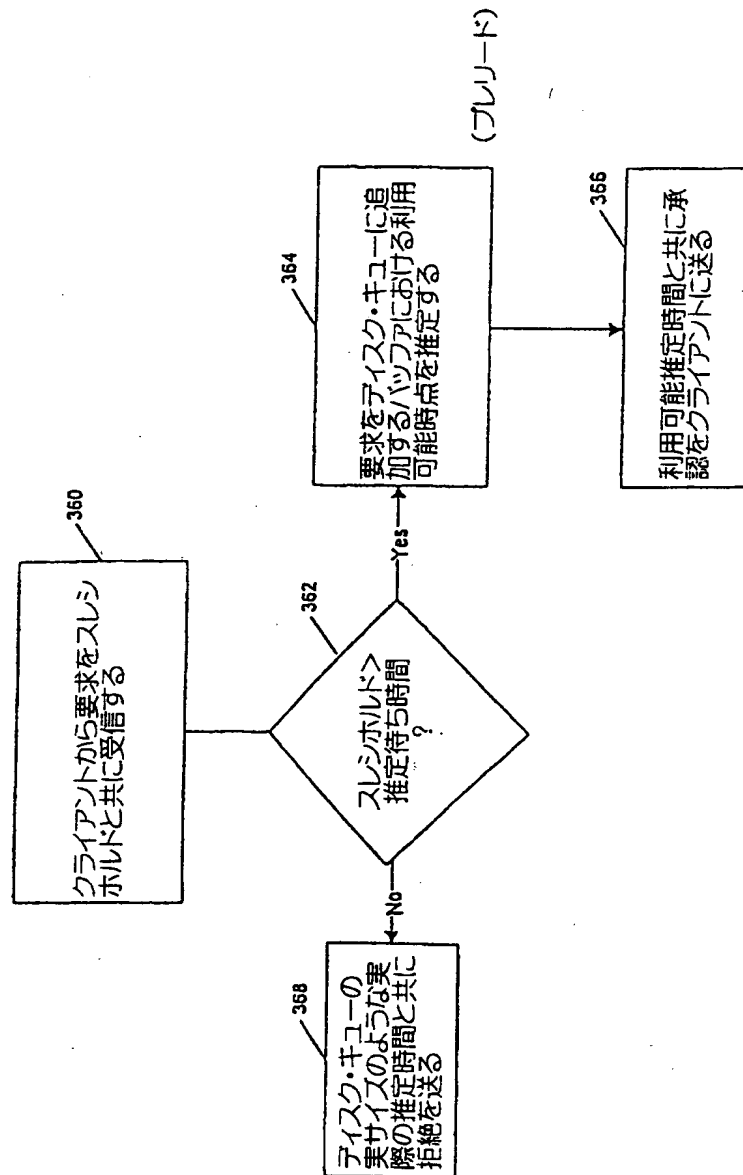
【補正対象項目名】図18



【補正方法】変更

【補正内容】

【図18】



【手続補正18】

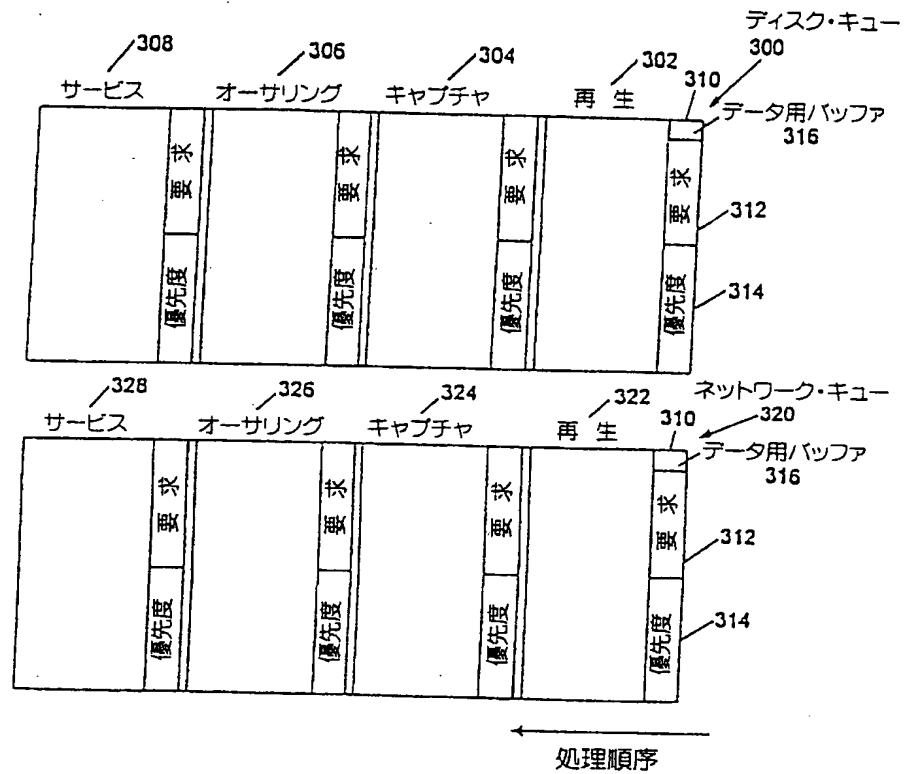
【補正対象書類名】図面

【補正対象項目名】図19

【補正方法】変更

【補正内容】

【図19】



【手続補正19】

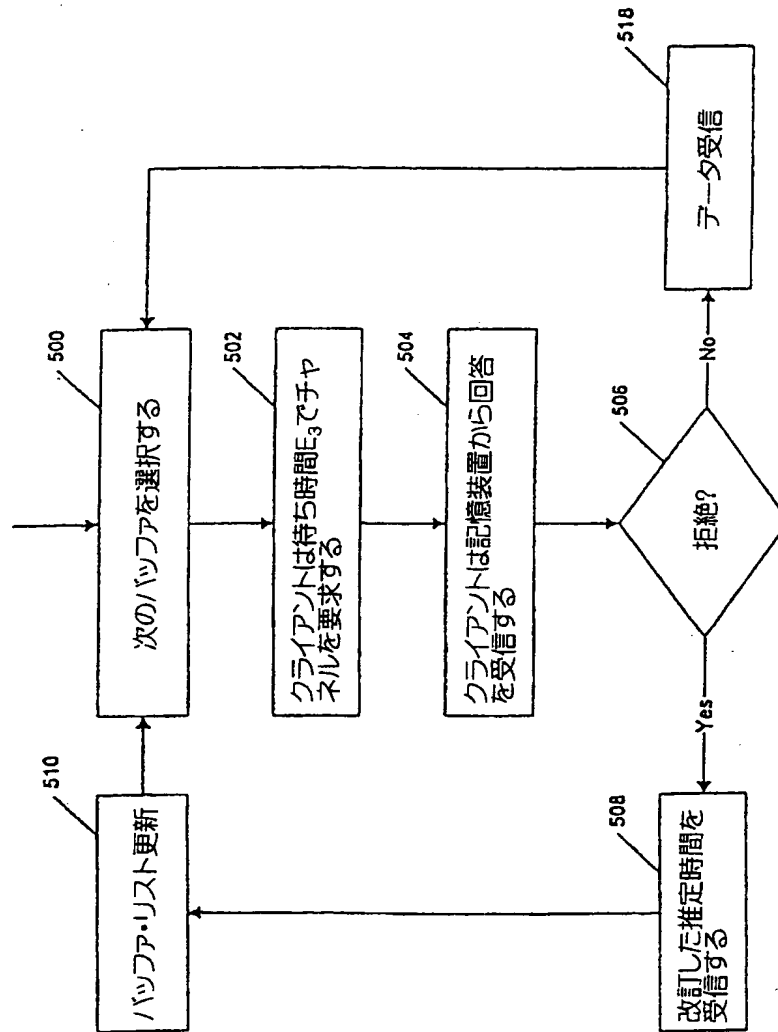
【補正対象書類名】図面

【補正対象項目名】図20

【補正方法】変更

【補正内容】

【図20】



【手続補正20】

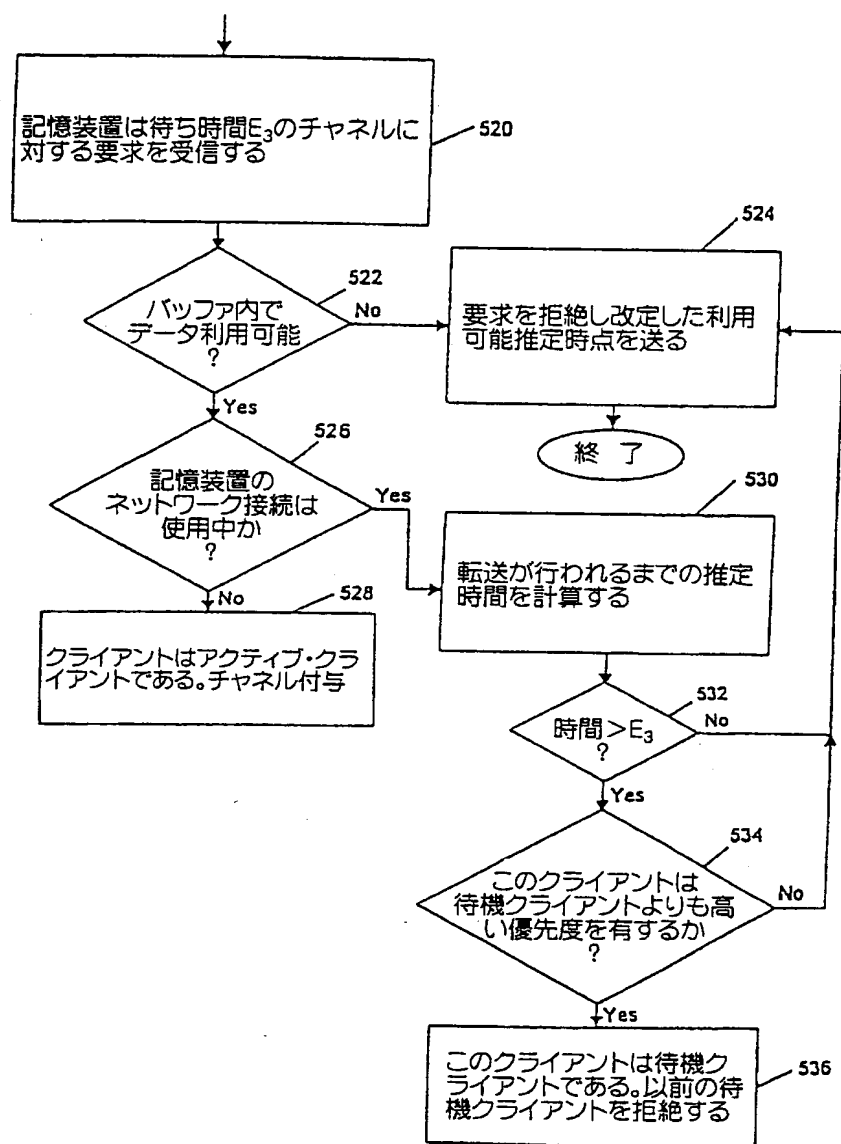
【補正対象書類名】図面

【補正対象項目名】図21

【補正方法】変更

【補正内容】

【図21】



【手続補正21】

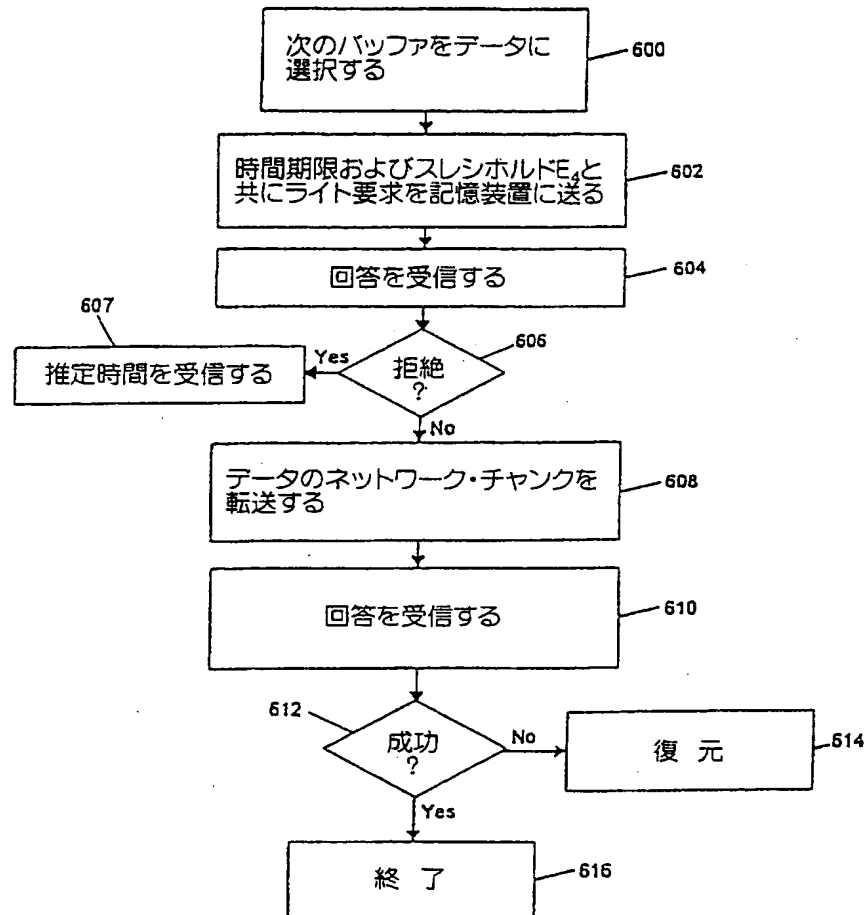
【補正対象書類名】図面

【補正対象項目名】図22

【補正方法】変更

【補正内容】

【図22】



【手続補正22】

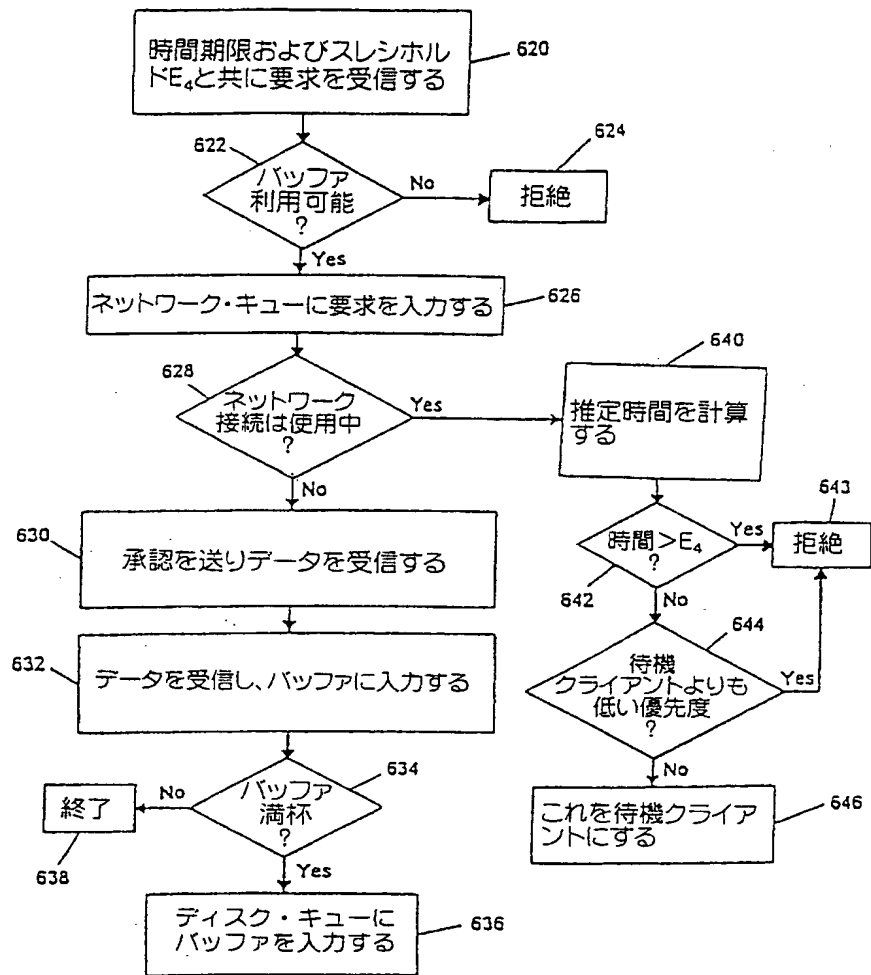
【補正対象書類名】図面

【補正対象項目名】図23

【補正方法】変更

【補正内容】

【図23】



【手続補正23】

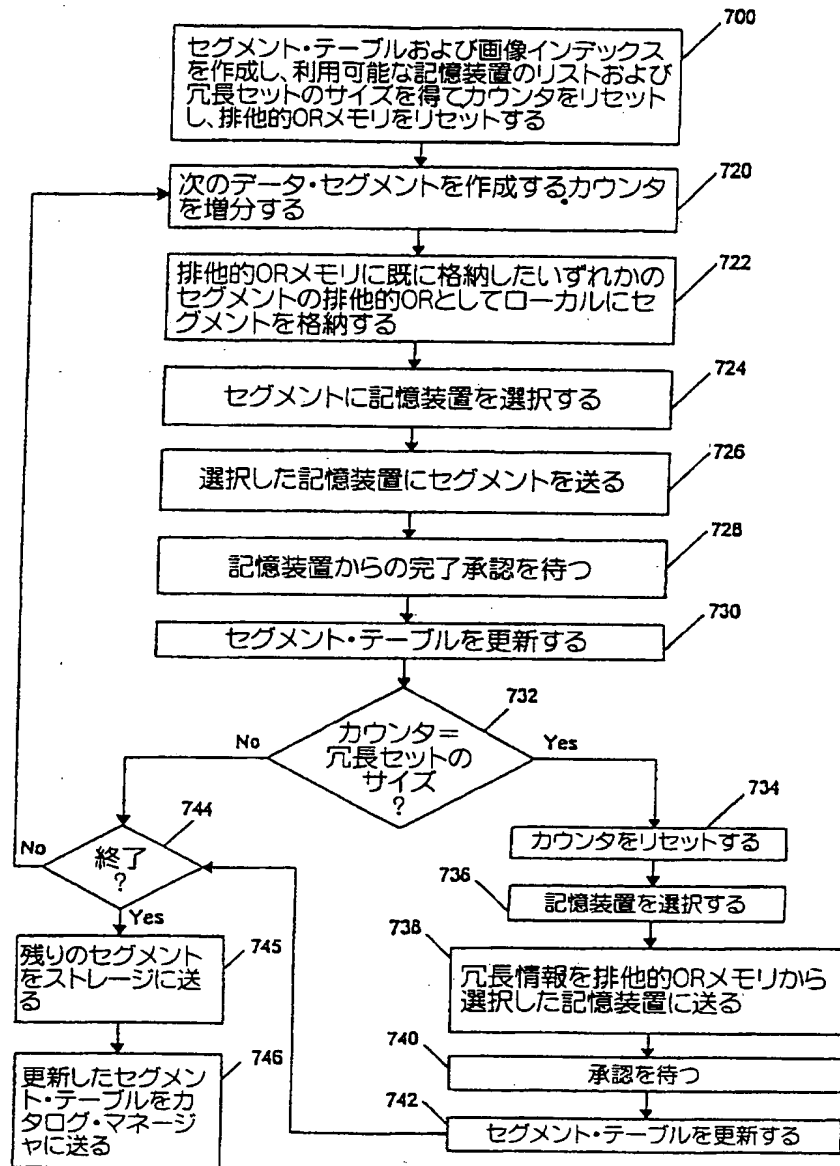
【補正対象書類名】図面

【補正対象項目名】図24

【補正方法】変更

【補正内容】

【図24】



【手続補正24】

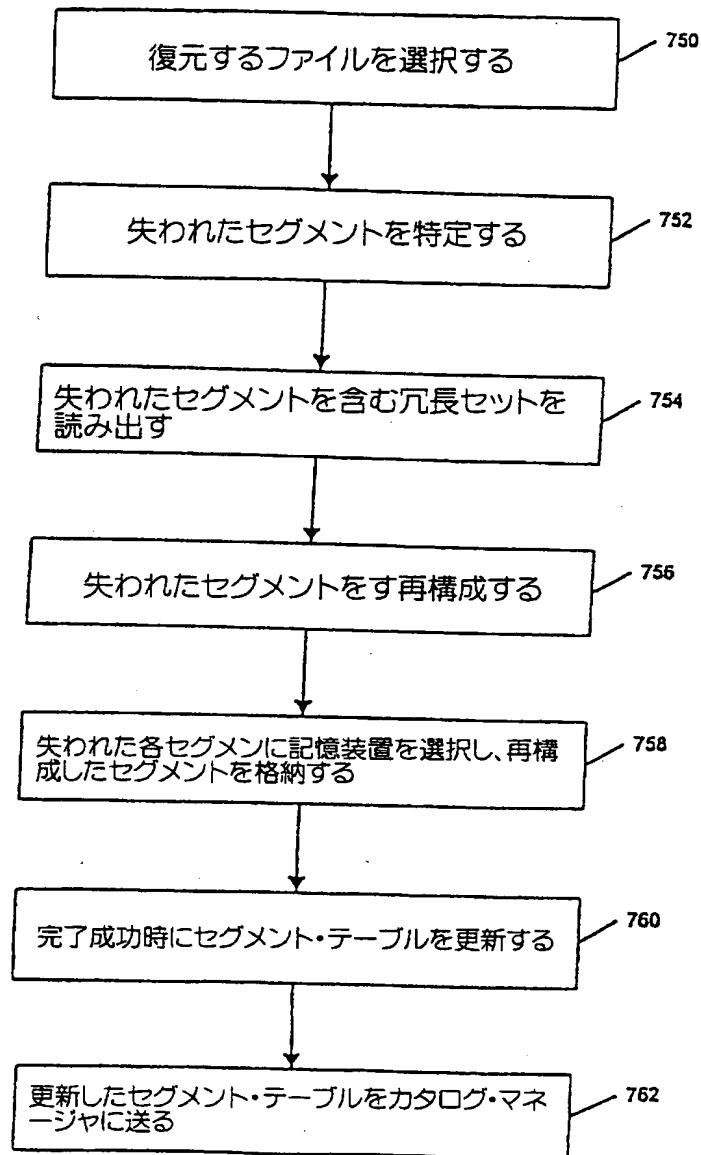
【補正対象書類名】図面

【補正対象項目名】図25

【補正方法】変更

【補正内容】

【図25】





## 【国際調査報告】

## INTERNATIONAL SEARCH REPORT

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> IPC 6 G06F11/20 G06F11/10 H04N7/173		International Application No. PCT/US 98/27199
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b> Minimum documentation searched (classification system followed by classification symbols) IPC 6 G06F H04N		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practical, search terms used)		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	R. TEWARI ET AL.: "High Availability in Clustered Multimedia Servers" PROC. 12TH INT. CONF. ON DATA ENGINEERING, 26 February 1996, pages 645-654, XP000632617 new orleans, la, usa see the whole document	1-10, 25-37
X	EP 0 701 198 A (STARLIGHT NETWORKS) 13 March 1996	1,5-9
Y	see column 7, line 43 - column 8, line 40 see column 27, line 18 - line 47 --- -/-	16
<input checked="" type="checkbox"/> Further documents are listed in the continuation of box C. <input checked="" type="checkbox"/> Patent family members are listed in annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier document but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art. "Z" document member of the same patent family		
Date of the actual completion of the international search 8 June 1999		Date of mailing of the international search report 23/06/1999
Name and mailing address of the ISA European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Tr. 31 851 apo nl, Fax (+31-70) 340-3018		Authorized officer Absalom, R

Form PCT/ISA-210 (second sheet) (July 1992)

## INTERNATIONAL SEARCH REPORT

International Application No.  
PCT/US 98/27199

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT		
Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	R. FLYNN ET AL.: "Disk Striping and Block Replication Algorithms for Video File Servers" PROC. INT. CONF. ON MULTIMEDIA COMPUTING AND SYSTEMS, 17 June 1996, pages 590-597, XP002105211 Hiroshima, Japan see the whole document	1-6
Y	EP 0 740 247 A (HEWLETT-PACKARD COMPANY) 30 October 1996	16
A	see abstract	10-15
A	GB 2 299 424 A (MITSUBISHI DENKI KABUSHIKI KAISHA) 2 October 1996 see the whole document	1-37
A	B. NARENDHAN: "Data Distribution Algorithms for Load Balanced Fault-Tolerant Web Access" PROC. 16TH SYMP. ON RELIABLE DISTRIBUTED SYSTEMS, 22 October 1997, pages 97-106, XP002105212 Durham, NC, USA see the whole document	1-37
A	S. GHANDEHARIZADEH ET AL.: "Continuous Retrieval of Multimedia Data Using Parallelism" IEEE TRANS. ON KNOWLEDGE AND DATA ENGINEERING, vol. 5, no. 4, August 1993, pages 658-669, XP002105213 USA see the whole document	1-37
A	US 5 559 764 A (CHEN ET AL.) 24 September 1996 cited in the application	

1

Form PCT/ISA/216 (continuation of second sheet) (July 1992)

## INTERNATIONAL SEARCH REPORT

Information on patent family members

International Application No.

PCT/US 98/27199

Patent document cited in search report		Publication date	Patent family member(s)		Publication date
EP 701198	A	13-03-1996	US	5732239 A	24-03-1998
EP 740247	A	30-10-1996	US	5592612 A	07-01-1997
			JP	9026854 A	28-01-1997
GB 2299424	A	02-10-1996	JP	8329021 A	13-12-1996
			US	5630007 A	13-05-1997
US 5559764	A	24-09-1996	EP	0697660 A	21-02-1996
			JP	8069360 A	12-03-1996

## フロントページの続き

(51)Int.Cl. <sup>7</sup>	識別記号	F I	テマコード(参考)
G 0 6 F 12/16	3 2 0	G 0 6 F 12/16	3 2 0 L 5 C 0 6 4
H 0 4 N 7/173	6 1 0	H 0 4 N 7/173	6 1 0 A

(31)優先権主張番号 09/054, 761

(32)優先日 平成10年4月3日(1998. 4. 3)

(33)優先権主張国 米国 (US)

(81)指定国 EP(AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), AU, CA, CN, DE, GB, JP

(71)出願人 Metropolitan Technology Park, One Park West, Tewksbury, Massachusetts 01876, United States of America

(72)発明者 ラビノウィッツ, スタンレー  
アメリカ合衆国マサチューセッツ州01886,  
ウエストフォード, ヴァイン・ブルック・  
ロード 12(72)発明者 ジェイコブス, ハーバート・アール  
アメリカ合衆国ニューハンプシャー州  
03051, ハドソン, サンライズ・ドライブ  
17(72)発明者 ジレット, リチャード・ベイカー, ジュニア  
アメリカ合衆国マサチューセッツ州01886,  
ウエストフォード, プレザベーション・ウ  
エイ 30(72)発明者 ファッシアノ, ピーター・ジェイ  
アメリカ合衆国マサチューセッツ州01760,  
ネイティック, コーチマン・レイン 30Fターム(参考) 5B001 AA00 AB01 AC01 AD03  
5B014 EA04 EB04 FB03 FB04 GC07  
GD23 GD32 GD33 HA09  
5B018 GA01 HA11 MA11 MA15  
5B065 BA01 EA03 EA12 EA19 EA24  
EA31 EA35 ZA08 ZA15  
5B082 DE05 HA01  
5C064 BA07 BB06 BC18 BD02 BD13

## 【要約の続き】

ての記憶装置上で均衡化する。この技法の組み合わせの結果、多数のアプリケーションおよび多数の記憶装置間で双方向にスケーラブルに多数の独立した高帯域データ・ストリームを転送可能なシステムが得られる。